



LUDWIG-  
MAXIMILIANS-  
UNIVERSITÄT  
MÜNCHEN

INSTITUT FÜR STATISTIK



Martin Slawski, Wolfgang zu Castell & Gerhard Tutz

# Feature Selection Guided by Structural Information

Technical Report Number 051, 2009  
Department of Statistics  
University of Munich

<http://www.stat.uni-muenchen.de>



# Feature Selection Guided by Structural Information

Martin Slawski<sup>1,2,3\*</sup>

Wolfgang zu Castell<sup>4†</sup>

Gerhard Tutz<sup>1‡</sup>

<sup>1</sup> Department of Statistics, Ludwig-Maximilians-University Munich, Germany

<sup>2</sup> Sylvia Lawry Centre, Munich, Germany

<sup>3</sup> Department of Computer Science, Saarland University, Saarbrücken, Germany

<sup>4</sup> Helmholtz Zentrum München, German Research Center for Environmental Health, Neuherberg, Germany

## Abstract

In generalized linear regression problems with an abundant number of features, lasso-type regularization which imposes an  $\ell^1$ -constraint on the regression coefficients has become a widely established technique. Crucial deficiencies of the lasso were unmasked when Zhou and Hastie (2005) introduced the elastic net. In this paper, we propose to extend the elastic net by admitting general nonnegative quadratic constraints as second form of regularization. The generalized ridge-type constraint will typically make use of the known association structure of features, e.g. by using temporal- or spatial closeness.

We study properties of the resulting 'structured elastic net' regression estimation procedure, including basic asymptotics and the issue of model selection consistency. In this vein, we provide an analog to the so-called 'irrepresentable condition' which holds for the lasso. An oracle property is established by incorporating a scaled  $\ell^1$ -constraint. Moreover, we outline algorithmic solutions for the structured elastic net within the generalized linear model family. The rationale and the performance of our approach is illustrated by means of simulated and real world data.

*keywords:* generalized linear model, regularization, sparsity,  $p \gg n$ , lasso, elastic net, random fields, consistency, epiconvergence, model selection, signal regression.

## 1 Introduction

We consider regression problems with a linear predictor. Let  $\mathbb{X} = (X_1, \dots, X_p)^\top$  be a random vector of real-valued features/predictors and let  $Y$  be a random response variable taking values in a set  $\mathcal{Y} \subseteq \mathbb{R}$ . Given a realization  $\mathbf{x} = (x_1, \dots, x_p)^\top$  of  $\mathbb{X}$ , a prediction  $\hat{y}$  for the response is obtained via a linear predictor

$$f(\mathbf{x}; \beta_0, \boldsymbol{\beta}) = \beta_0 + \mathbf{x}^\top \boldsymbol{\beta}, \quad \boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^\top,$$

and a function  $\zeta : \mathbb{R} \rightarrow \mathcal{Y}$  such that  $\hat{y} = \zeta(f(\mathbf{x}))$ . Given an i.i.d. sample  $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$  from  $(\mathbb{R}^p \times \mathcal{Y})^n$ , an optimal set of coefficients  $\hat{\beta}_0, \hat{\boldsymbol{\beta}} = (\hat{\beta}_1, \dots, \hat{\beta}_p)^\top$  can be determined by minimization of a criterion of the form

$$(\hat{\beta}_0, \hat{\boldsymbol{\beta}}) = \underset{(\beta_0, \boldsymbol{\beta})}{\operatorname{argmin}} \sum_{i=1}^n L(y_i, f(\mathbf{x}_i; \beta_0, \boldsymbol{\beta})), \quad (1.1)$$

---

\*[ms@cs.uni-sb.de](mailto:ms@cs.uni-sb.de)

†[castell@helmholtz-muenchen.de](mailto:castell@helmholtz-muenchen.de)

‡[tutz@stat.uni-muenchen.de](mailto:tutz@stat.uni-muenchen.de)

where  $L : \mathcal{Y} \times \mathbb{R} \rightarrow \mathbb{R}_0^+$  is a loss function, assumed to be continuous and convex in the second argument. The loss function is chosen according to the specific prediction problem, so that large loss represents bad fit to the observed sample  $S$ . Approach (1.1) usually yields poor estimates  $\hat{\beta}_0, \hat{\beta}$  if  $n$  is not one order of magnitude larger than  $p$ . In particular, if  $p \gg n$ , approach (1.1) is not well-defined in the sense that there exist infinitely many minimizers  $\hat{\beta}_0, \hat{\beta}$ . One way to cope with a small  $n/p$  ratio is to employ a regularizer  $\Omega(\beta)$ . A traditional approach due to Hoerl and Kennard (1970) minimizes the loss in Eq. (1.1) subject to an  $\ell^2$ -constraint on  $\beta$ . In the situation that  $\beta$  is supposed to be sparse, Tibshirani (1996) proposed, under the acronym 'lasso', to work with an  $\ell^1$ -constraint, i.e. one maximizes the loss subject to  $\Omega(\beta) = \|\beta\|_1 < s$ ,  $s > 0$ . The latter is particularly attractive if one is interested in feature selection, since one obtains estimates  $\hat{\beta}_j$ ,  $j \in \{1, \dots, p\}$ , which equal exactly zero, such that feature  $j$  does not contribute to prediction, for which we say that feature  $j$  is 'not selected'. Continuous shrinkage (Fan and Li (2001)) and the existence of efficient algorithms (Efron et al. (2004), Genkin et al. (2007)) for determining the coefficients are further virtues of the lasso. Its limitations have recently been revealed by several researchers. Zhou and Hastie (2005) pointed out that the lasso degenerates in the  $p \gg n$  setting, where the lasso is able to select at most  $n$  features (Rosset et al. (2004)). Furthermore, Zhou and Hastie stated that the lasso does not distinguish between 'irrelevant' and 'relevant but redundant' features. In particular, if there is a group of correlated features, then the lasso tends to select one arbitrary member of the group while ignoring the remainder. The combined regularizer of the elastic net  $\Omega(\beta) = \alpha \|\beta\|_1 + (1 - \alpha) \|\beta\|^2$ ,  $\alpha \in (0, 1)$  is shown to provide remedy in this regard. A second double regularizer - tailored to one-dimensional signal regression - is employed by the fused lasso (Tibshirani et al. (2005)), who propagate  $\Omega(\beta) = \alpha \|\beta\|_1 + (1 - \alpha) \|\mathbf{D}\beta\|_1$ , where

$$\begin{aligned} \mathbf{D} : \quad \mathbb{R}^p &\rightarrow \mathbb{R}^{p-1} \\ (\beta_1, \dots, \beta_p)^\top &\mapsto ([\beta_2 - \beta_1], \dots, [\beta_p - \beta_{p-1}])^\top \end{aligned} \quad (1.2)$$

is the first forward difference operator. The total variation regularizer is meaningful whenever there is an order relation, notably a temporal one, among the features. The fused lasso has a property which can be beneficial for interpretation: it automatically clusters the features, since the sequence  $\hat{\beta}_1, \dots, \hat{\beta}_p$  is blockwise constant.

In this paper, we study a regularizer which is intermediate between the elastic net and the fused lasso. Our regularizer combines an  $\ell^1$ -constraint with a quadratic form:

$$\Omega(\beta) = \alpha \|\beta\|_1 + (1 - \alpha) \beta^\top \mathbf{\Lambda} \beta, \quad (1.3)$$

where  $\mathbf{\Lambda} = (l_{jj'})_{1 \leq j, j' \leq p}$  is assumed to be symmetric and positive semidefinite. Setting  $\mathbf{\Lambda} = \mathbf{I}$  yields the elastic net. Therefore, expression (1.3) will be referred to as structured elastic net regularizer. The inclusion of  $\mathbf{\Lambda}$  aims at capturing the a priori association structure (if available) of the features in more generality than the fused lasso. The structured elastic net estimator is defined as

$$(\hat{\beta}_0, \hat{\beta}) = \underset{(\beta_0, \beta)}{\operatorname{argmin}} \sum_{i=1}^n L(y_i, f(\mathbf{x}_i; \beta_0, \beta)) \quad (1.4)$$

$$\text{subject to } \alpha \|\beta\|_1 + (1 - \alpha) \beta^\top \mathbf{\Lambda} \beta \leq s, \quad \alpha \in (0, 1), \quad s > 0,$$

which is equivalent to the Lagrangian formulation

$$(\hat{\beta}_0, \hat{\beta}) = \underset{(\beta_0, \beta)}{\operatorname{argmin}} \sum_{i=1}^n L(y_i, f(\mathbf{x}_i; \beta_0, \beta)) + \lambda_1 \|\beta\|_1 + \lambda_2 \beta^\top \mathbf{\Lambda} \beta, \quad \lambda_1, \lambda_2 > 0. \quad (1.5)$$

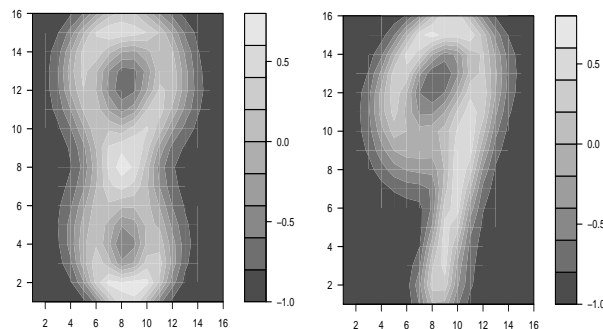
The rest of the paper is organized as follows: in Section 2, we discuss the choice of the matrix  $\mathbf{\Lambda}$ , followed by an analysis of some important properties of our proposal (1.5) in Section 3. Section 4 is devoted to asymptotics and consistency questions, motivating the introduction of the *adaptive* structured elastic net. Section 5 presents three different algorithms for computing the minimizers (1.5) in the generalized linear model family and addresses fundamental questions of model selection and -inference. The practical performance of the structured elastic net is contained in Section 6. Section 7 concludes with a brief discussion and outlook. All proofs can be found in the Appendix.

## 2 Structured features

### 2.1 Motivation

A considerable fraction of contemporary regression problems is characterized by a large number of features, which is either of the same order of magnitude as the sample size or even several orders larger ( $p \gg n$ ). Common instances thereof are feature sets consisting of sampled signals, pixels of an image, spatially sampled data, or gene expression intensities. Beside high dimensionality of the feature space, these examples have in common that the feature set can be arranged according to an a priori association structure. If a sampled signal does not vary rapidly, the influence of nearby sampling points on the response can be expected to be similar; correspondingly, this applies to adjacent pixels of an image, or, more general, to any other form of spatially linked features. In genomics, genes can be categorized into functional groups, or one has prior knowledge of their functions and interactions within biochemical reaction chains, so called pathways.

Figures 1 and 2 display two well-known examples, phoneme- and handwritten digit classification. These examples are well-appt to illustrate the idea of the structured elastic net regularizer, since it is sensible to assume that the prediction problem is not only characterized by smoothness with respect to a given structure, but also by sparsity: in the phoneme classification example, visually only the first hundred frequencies seem to carry information relevant to the prediction problem. A similar rationale applies to the second example, where the arc of the numeral eight in the lower half of the picture is the eminent characteristic that admits a distinction from the numeral nine.



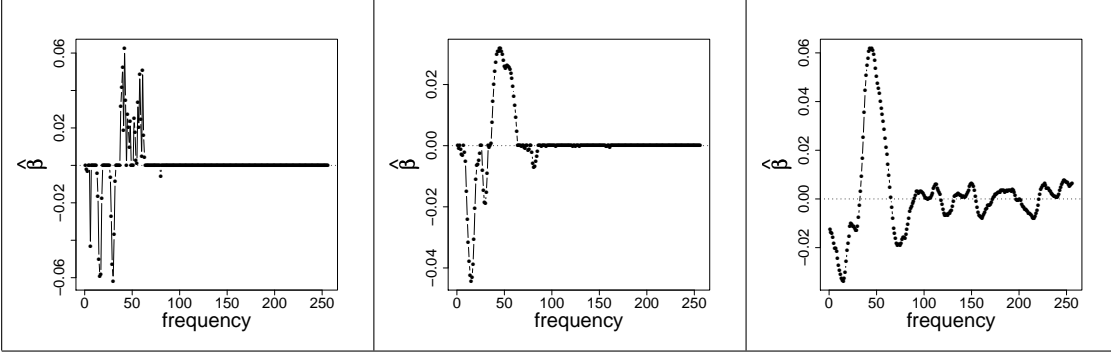


Figure 1: Phoneme data (Hastie et al. (1995)). The upper panel shows several thousand log-periodograms of the speech frames for the phonemes 'aa' (as occurring in 'dark') and 'ao' (as occurring in 'water'). The classwise means are given by thick lines. We use linear logistic regression to predict the phoneme given a log-periodogram. The lower panel depicts the resulting coefficients when using the lasso (left panel), a first-order difference penalty (right panel), and a combination thereof, which we term 'structured elastic net' (middle panel).

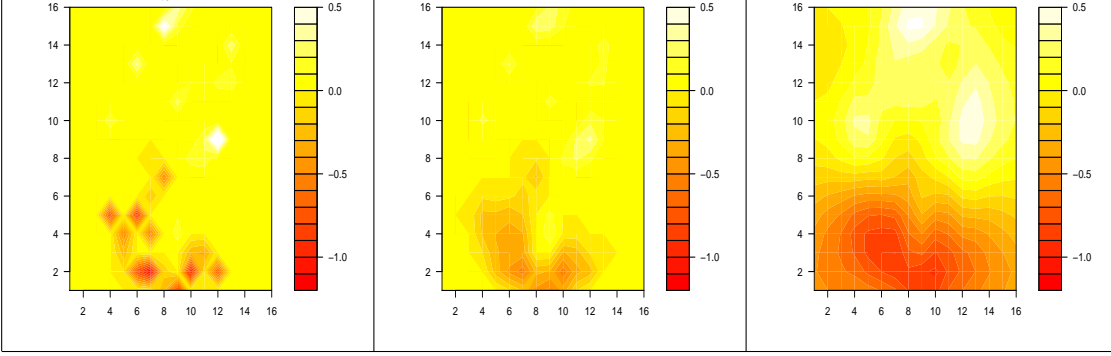


Figure 2: Handwritten digit recognition dataset (Le Cun et al. (1989)). One observation is given by a greyscale image composed of  $16 \times 16$  pixels. The upper panel shows the contour of the pixel-wise means for the numerals '8' and '9'. We use a training set of 1500 observations of eights and nines as input for linear logistic regression. The lower panel depicts the coefficient surfaces for the lasso (left panel), a discrete Laplacian penalty according to the grid structure (right panel), and a combination, the structured elastic net (middle panel).

## 2.2 Gauss-Markov random fields

Given a large, but structured set of features, its structure can be exploited to cope with high dimensionality in regression estimation. The estimands  $\{\beta_j\}_{j=1}^p$  form a finite set such that their prior dependence structure can conveniently be described by means of a graph  $\mathcal{G} = (V, E)$ ,  $V = \{\beta_1, \dots, \beta_p\}$ ,  $E \subset V \times V$ . We exclude loops, i.e.  $(\beta_j, \beta_j) \notin E$  for all  $j$ . The edges may additionally be weighted by a function  $w : E \rightarrow \mathbb{R}$ ,  $w((\beta_j, \beta_{j'})) = w((\beta_{j'}, \beta_j))$  for all edges in  $E$ . We will use the notation  $\beta_j \sim \beta_{j'}$  to express that  $\beta_j$  and  $\beta_{j'}$  are connected by an edge in  $\mathcal{G}$ . The weight function can be extended to a function on  $V \times V$  by setting  $w((\beta_j, \beta_{j'})) = w((\beta_{j'}, \beta_j)) = 0$  if  $(\beta_j, \beta_{j'}) \notin E$ .

The graph is interpreted in terms of Gauss-Markov random fields (Besag (1974); Rue and

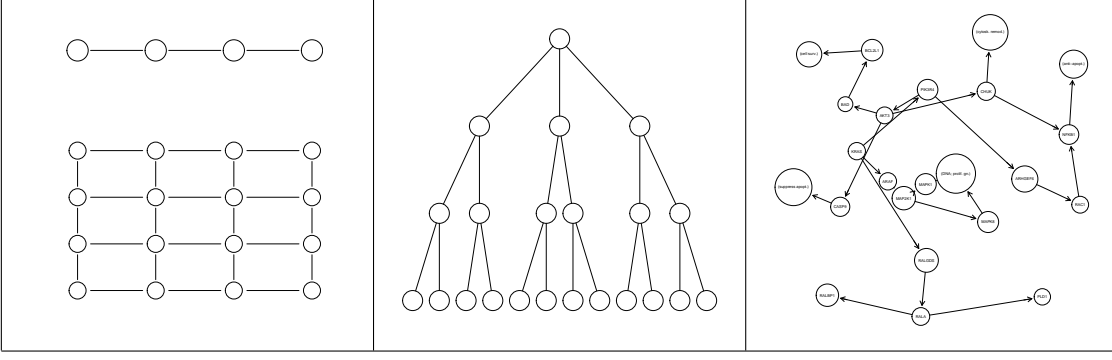


Figure 3: A collection of some graphs. A path and a grid (left panel), a rooted tree (middle panel) and an irregular graph describing a part of the so-called MAPK signaling pathway (right panel).

Held (2001)). In our setup, the pairwise Markov property reads

$$\neg \beta_j \sim \beta_{j'} \Leftrightarrow \beta_j \perp\!\!\!\perp \beta_{j'} \mid V \setminus \{\beta_j, \beta_{j'}\}, \quad (2.1)$$

with  $\perp\!\!\!\perp$  denoting conditional independence. Property (2.1) is conform to the following choice for the precision matrix  $\mathbf{\Lambda} = (l_{jj'})_{1 \leq j, j' \leq p}$ :

$$l_{jj'} = \begin{cases} \sum_{k=1}^p |w((\beta_j, \beta_k))| & \text{if } j = j', \\ -w((\beta_j, \beta_{j'})) & \text{if } j \neq j', \end{cases} \quad (2.2)$$

which is singular in general. If  $\text{sign}\{w((\beta_j, \beta_{j'}))\} \geq 0$  for all  $(\beta_j, \beta_{j'})$  in  $E$ , then  $\mathbf{\Lambda}$  as given in Eq. (2.2) is known as the combinatorial graph Laplacian in spectral graph theory (Chung (1997)). It is straightforward to verify the following properties.

- $$\beta^\top \mathbf{\Lambda} \beta = \sum_{\beta_j \sim \beta_{j'}} |w(\beta_j, \beta_{j'})| (\beta_j - \text{sign}\{w((\beta_j, \beta_{j'}))\} \beta_{j'})^2 \geq 0, \quad (2.3)$$

where the sum is over all distinct edges in  $\mathcal{G}$ , and 'distinct' is understood with respect to the relation  $(\beta_j, \beta_{j'}) = (\beta_{j'}, \beta_j)$  for all  $j, j'$ .

- If  $\mathcal{G}$  is connected and  $\text{sign}\{w((\beta_j, \beta_{j'}))\} \geq 0$  for all  $(\beta_j, \beta_{j'})$  in  $E$ , the nullspace of  $\mathbf{\Lambda}$  is spanned by the vector of ones  $\mathbf{1}$ .

While we have started in full generality, the choice  $w((\beta_j, \beta_{j'})) \in \{0, 1\}$  for all  $j, j'$  will frequently be the standard choice in practice. In this case, the quadratic form captures local fluctuations of  $\beta$  w.r.t.  $\mathcal{G}$ . As a simple example, one may take  $\mathcal{G}$  as the path on  $p$  vertices so that expression (2.3) equals the summed squared forward differences

$$\sum_{j=2}^p (\beta_j - \beta_{j-1})^2 = \|\mathbf{D}\beta\|^2 = \beta^\top \mathbf{D}^\top \mathbf{D} \beta, \quad (2.4)$$

where  $\mathbf{D}$  is defined in Eq. (1.2). More complex graphical structures can be generated from simple ones using the notion of Cartesian products of graphs. Given two unweighted

graphs  $\mathcal{G} = (V, E)$  and  $\mathcal{G}' = (V', E')$ , their Cartesian product is defined as

$$\begin{aligned}\mathcal{G} \times \mathcal{G}' &= (V_{\times}, E_{\times}), \quad V_{\times} = V \times V', \\ E_{\times} &= \{((j, k), (j', k')) \in V_{\times} \times V_{\times} : (j, j') \in E \wedge k = k' \vee (k, k') \in E' \wedge j = j'\}.\end{aligned}\tag{2.5}$$

In terms of adjacency matrices  $\mathbf{A} = (a_{jj'} = I(j \overset{\mathcal{G}}{\sim} j'))$  and  $\mathbf{A}' = (a_{kk'} = I(k \overset{\mathcal{G}'}{\sim} k'))$ , where  $I$  is the indicator function, definition (2.5) can more comfortably be expressed as

$$\mathbf{A}_{\times} = \mathbf{A} \otimes \mathbf{I}' + \mathbf{I} \otimes \mathbf{A}',\tag{2.6}$$

where  $\mathbf{A}_{\times}$  denotes the adjacency matrix of  $\mathcal{G} \times \mathcal{G}'$ ,  $\mathbf{I}$  and  $\mathbf{I}'$  are identity matrices with dimensions equal to the number of vertices in  $\mathcal{G}$  and  $\mathcal{G}'$ , respectively, and  $\otimes$  denotes the Kronecker product. It is easy to see that this construction yields a regular grid on  $p \times p'$  vertices if  $\mathcal{G}$  is chosen as  $p$ -path and  $\mathcal{G}'$  as  $p'$ -path, in which case  $\mathbf{\Lambda}$  equals the usual discretization of the Laplacian  $\Delta$  acting on functions defined on  $\mathbb{R}^2$ . Regularizers built up from discrete differences have already seen frequent use in high-dimensional regression estimation. Examples comprise penalized discriminant analysis (Hastie et al. (1995)) and spline smoothing (Eilers and Marx (1996)).

### 3 Properties

#### 3.1 Bayesian and geometric interpretation

In the setup of Section 1, consider the regularizer

$$\Omega(\boldsymbol{\beta}) = \lambda_1 \|\boldsymbol{\beta}\|_1 + \lambda_2 \boldsymbol{\beta}^{\top} \mathbf{\Lambda} \boldsymbol{\beta}, \quad \lambda_1, \lambda_2 > 0.$$

It has a nice Bayesian interpretation when the loss function  $L$  is of the form

$$L(y, f(\mathbf{x}; \beta_0, \boldsymbol{\beta})) = \phi^{-1}(b(f(\mathbf{x})) - yf(\mathbf{x})) + c(y, \phi),\tag{3.1}$$

i.e. the loss function equals the negative log-likelihood of a generalized linear model in canonical parametrization, which will primarily be studied in this paper. Models of this class are characterized by

$$\begin{aligned}Y|\mathbb{X} = \mathbf{x} &\sim \text{simple exponential family}, \\ \hat{y} = \mathbb{E}[Y|\mathbb{X} = \mathbf{x}] &= \mu = \frac{d}{df}b(f(\mathbf{x})), \\ \text{var}[Y|\mathbb{X} = \mathbf{x}] &= \phi \frac{d^2}{df^2}b(f(\mathbf{x})).\end{aligned}\tag{3.2}$$

The form (3.1) is versatile, including classical linear regression with Gaussian errors, logistic regression for classification, and Poisson regression for count data. Given a loss from the class (3.1), the regularizer  $\Omega(\boldsymbol{\beta})$  can be interpreted as combined Laplace (double exponential)-Gaussian prior  $p(\boldsymbol{\beta}) \propto \exp(-\Omega(\boldsymbol{\beta}))$ , for which the structured elastic net estimator (1.5), provided  $p(\beta_0) \propto 1$ , is the maximum posterior (MAP) estimator given the sample  $S$ . Note that if  $p > n$ , depending on the matrix  $\mathbf{\Lambda}$ ,  $p(\boldsymbol{\beta})$  is not necessarily proper. It is instructive to consider two predictors, i.e.  $\boldsymbol{\beta} = (\beta_1, \beta_2)^{\top}$ . Figure 4 gives a geometric interpretation for the basic choices  $\mathbf{\Lambda} = \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix}$  and  $\mathbf{\Lambda} = \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}$ , corresponding

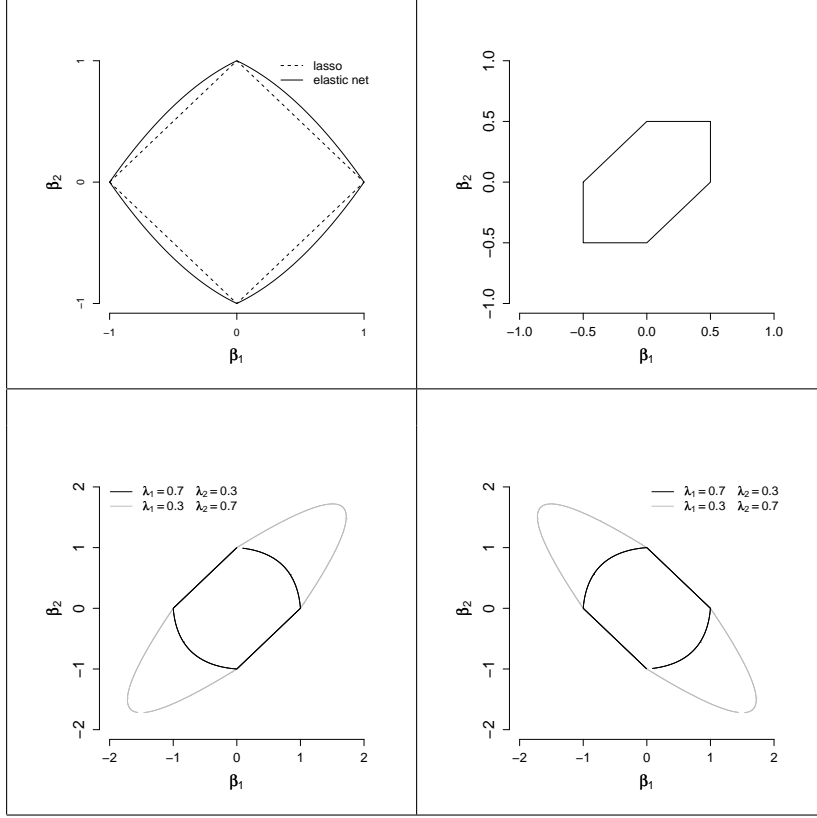


Figure 4: Geometry of the lasso and the elastic net (upper left panel), the fused lasso (upper right panel) and the structured elastic net (positive prior correlation: left panel, negative prior correlation: right panel).

to positive- and negative prior correlation, respectively.

The contour lines of the structured elastic net penalty contain elements of a diamond and an ellipsoid. The higher  $\lambda_2$  in relation to  $\lambda_1$ , the ellipsoidal part becomes more narrower and more stretched. The sign of the off-diagonal element of  $\mathbf{\Lambda}$  determines the orientation of the ellipsoidal part.

### 3.2 Grouping properties

For the elastic net, [Zhou and Hastie \(2005\)](#) provided an upper bound on the absolute distances  $|\hat{\beta}_j^{\text{elastic net}} - \hat{\beta}_{j'}^{\text{elastic net}}|$ ,  $j, j' = 1, \dots, p$ , in terms of the sample correlations, to which Zhou and Hastie referred to as 'grouping property'. We provide similar bounds here. For what follows, let  $S$  be a sample as in Section 1. We introduce a design matrix  $\mathbf{X} = (x_{ij})_{1 \leq i \leq n, 1 \leq j \leq p}$  and denote by  $\mathbf{X}_j = (x_{1j}, \dots, x_{nj})^\top$  the realizations of predictor  $j$  in

$S$ , and the response vector is defined by  $\mathbf{y} = (y_1, \dots, y_n)^\top$ . For the remainder of this section, we assume that the responses are centered and that the predictors are centered and standardized to unit Euclidean length w.r.t the sample  $S$ , i.e.

$$\sum_{i=1}^n y_i = \sum_{i=1}^n x_{ij} = 0, \quad \sum_{i=1}^n x_{ij}^2 = 1, \quad j = 1, \dots, p. \quad (3.3)$$

**Proposition 1.** *Let  $p = 2$ , let the loss function be of the form (3.1), let  $\rho = \langle \mathbf{X}_1, \mathbf{X}_2 \rangle$*



denote the sample correlation of  $\mathbf{X}_1$  and  $\mathbf{X}_2$ , and let  $\mathbf{\Lambda} = \frac{1}{2} \begin{pmatrix} 1 & s \\ s & 1 \end{pmatrix}$ ,  $s \in \{-1, 1\}$ . If  $-s\hat{\beta}_1\hat{\beta}_2 > 0$ , then

$$|\hat{\beta}_1 + s\hat{\beta}_2| \leq \frac{1}{2\lambda_2} \sqrt{2(1+s\rho)} \|\mathbf{y}\|.$$

In particular, Proposition 1 implies that if  $\mathbf{X}_1 = -s\mathbf{X}_2$ ,  $\hat{\beta}_1 = -s\hat{\beta}_2$ .

In a similar way as done in Wang et al. (2006), we can obtain a bound for loss functions which are either uniformly Lipschitz as function of the margin (classification) or as function of the residual (regression).

**Proposition 2.** *In the setup of Proposition 1 and for a positive finite constant  $C$ , assume that one of the following conditions holds:*

- (i)  $L(y, f(\mathbf{x})) = L(m)$ ,  $m = yf(\mathbf{x})$ ,  $y \in \{-1, 1\}$ ,  $|L(m) - L(m')| \leq C|m - m'|$  for all  $m, m'$ .
- (ii)  $L(y, f(\mathbf{x})) = L(r)$ ,  $r = y - f(\mathbf{x})$ ,  $|L(r) - L(r')| \leq C|r - r'|$  for all  $r, r'$ .

Then:

$$|\hat{\beta}_1 + s\hat{\beta}_2| \leq \frac{C}{\lambda_2} \sqrt{n2(1+s\rho)}.$$

Proposition 2 provides a grouping property for important loss functions, covering the hinge loss of support vector classification as well as least absolute deviation- and quantile regression.

### 3.3 Decorrelation

Let us now consider the important special case

$$L(y, f(\mathbf{x}; \boldsymbol{\beta})) = (y - \mathbf{x}^\top \boldsymbol{\beta})^2,$$

which corresponds to classical linear regression. The constant term  $\beta_0$  is omitted, since we work with centered data. The structured elastic net estimator can then be written as

$$\begin{aligned} \hat{\boldsymbol{\beta}} &= \underset{\boldsymbol{\beta}}{\operatorname{argmin}} -2\mathbf{y}^\top \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\beta}^\top [\mathbf{C} + \lambda_2 \mathbf{\Lambda}] \boldsymbol{\beta} + \lambda_1 \|\boldsymbol{\beta}\|_1, \quad \mathbf{C} = \mathbf{X}^\top \mathbf{X}, \\ &= \underset{\boldsymbol{\beta}}{\operatorname{argmin}} -2\mathbf{y}^\top \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\beta}^\top \tilde{\mathbf{C}} \boldsymbol{\beta} + \lambda_1 \|\boldsymbol{\beta}\|_1, \quad \tilde{\mathbf{C}} = \mathbf{X}^\top \mathbf{X} + \lambda_2 \mathbf{\Lambda}. \end{aligned} \tag{3.4}$$

Note that for standardized predictors,  $\mathbf{C}$  equals the matrix of sample correlations  $\rho_{jj'} = \langle \mathbf{X}_j, \mathbf{X}_{j'} \rangle$ ,  $j, j' = 1, \dots, p$ . With a large number of predictors or elements  $\rho_{jj'}$  with large  $|\rho_{jj'}|$ ,  $\mathbf{C}$  is known to yield severely unstable ordinary least squares (ols) estimates  $\hat{\beta}_j^{\text{ols}}$ ,  $j = 1, \dots, p$ . If the two underlying random variables  $X_j$  and  $X_{j'}$  are highly positively correlated, this will likely translate to high sample correlations of  $\mathbf{X}_j$  and  $\mathbf{X}_{j'}$ , which in turn yields a strongly negative correlation between  $\hat{\beta}_j^{\text{ols}}$  and  $\hat{\beta}_{j'}^{\text{ols}}$  and as consequence high variances  $\operatorname{var}[\hat{\beta}_j^{\text{ols}}]$  and  $\operatorname{var}[\hat{\beta}_{j'}^{\text{ols}}]$ . The modified matrix  $\tilde{\mathbf{C}}$  can be written as  $\tilde{\mathbf{C}} = \mathbf{V}_\Lambda^{1/2} \mathbf{R}_\Lambda \mathbf{V}_\Lambda^{1/2}$ ,  $\mathbf{V}_\Lambda = \operatorname{diag}(1 + \lambda_2 \sum_{k=1}^p |l_{1k}|, \dots, 1 + \lambda_2 \sum_{k=1}^p |l_{pk}|)$ , and the modified correlation matrix  $\mathbf{R}_\Lambda$  has entries

$$\rho_{\Lambda, jj'} = \frac{\rho_{jj'} + \lambda_2 l_{jj'}}{\sqrt{1 + \sum_{k=1}^p |l_{jk}|} \sqrt{1 + \sum_{k=1}^p |l_{j'k}|}}, \quad j, j' = 1, \dots, p.$$

In the light of Section 2, the entries of  $\mathbf{R}_\Lambda$  combine sample- and prior correlations. Decorrelation occurs if  $\rho_{jj'} \approx -\lambda_2 l_{jj'}$ , i.e. if prior- and sample correlations are in accordance. The grouping- and decorrelation effect of our proposal are visualized in Figure 5 for two predictors, in which case the minimizer (3.4) can be determined analytically (see Appendix B). The figure unmasks the weakness of ordinary least squares in the case of high sample correlations, as well as the tendency of the lasso not to jointly include the two predictors.

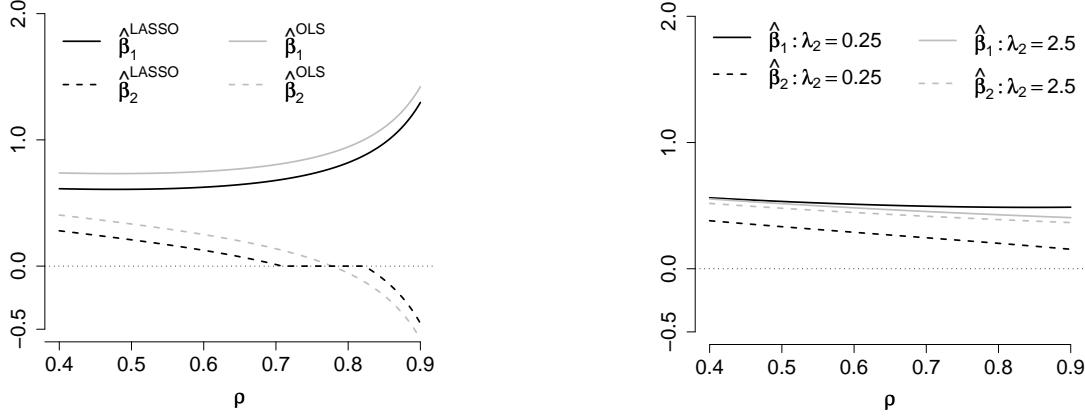


Figure 5: The case of two predictors for ordinary least squares, lasso and the structured elastic net. For both panels, the predictor-response correlations are set to  $\mathbf{X}_1^\top \mathbf{y} = 0.9$  and  $\mathbf{X}_2^\top \mathbf{y} = 0.7$ , respectively,  $\|\mathbf{y}\| = 1$ . The correlation of the predictors  $\rho = \mathbf{X}_1^\top \mathbf{X}_2$  varies from 0.4 to 0.9 (horizontal axis). The restriction  $\rho \in [0.4; 0.9]$  guarantees a valid correlation structure between  $\mathbf{X}_1, \mathbf{X}_2$  and  $\mathbf{y}$ . The left panel displays the behaviour of the ols estimator (grey) with  $\hat{\beta}_1^{\text{ols}}, \hat{\beta}_2^{\text{ols}}$  becoming more and more divergent, as does the lasso (black). The right panel depicts the same situation for the structured elastic net with  $\Lambda = (l_{jj'})_{1 \leq j, j' \leq 2}$ ,  $l_{11} = l_{22} = 1$ ,  $l_{12} = l_{21} = -1$  for  $\lambda_2 \in \{0.25, 2.5\}$  and  $\lambda_1$  as for the left panel. Note that the right panel does not contradict Proposition 1, since  $\rho \rightarrow 1$  implies  $\mathbf{X}_1^\top \mathbf{y} \rightarrow \mathbf{X}_2^\top \mathbf{y}$ , whereas  $\mathbf{X}_1^\top \mathbf{y}$  and  $\mathbf{X}_2^\top \mathbf{y}$  are kept fixed here. Further note that the situation displayed in the figure does not operate in terms of true coefficients, but studies instead the presence/absence of a grouping effect in dependence of a given data situation.

### 3.4 Uniqueness

Let us maintain the setup of the previous subsection. As pointed out by several authors, the lasso lacks uniqueness if  $p \gg n$  or if there is collinearity among  $\mathbf{X}_1, \dots, \mathbf{X}_p$ . In this respect, the lasso differs from the elastic net, which is always determined uniquely. Concerning our proposal, whether the minimizer (3.4) is unique for  $p \gg n$  depends on  $\Lambda$ . Roughly speaking, if  $\Lambda$  provides a sufficient amount of prior information, then we have a unique minimizer. Uniqueness typically fails when  $\Lambda$  has nonzero entries only for a small fraction of the predictors. A sufficient condition can be derived by the following consideration: since  $\Lambda$  is assumed to be symmetric positive semidefinite, there exists a factorization  $\Lambda = \mathbf{Q}^\top \mathbf{Q}$ . The stabilized covariance matrix in Eq. (3.4) can be written as

$$\tilde{\mathbf{C}} = \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}}, \quad \tilde{\mathbf{X}} = \begin{pmatrix} \mathbf{X} \\ \lambda_2^{1/2} \mathbf{Q} \end{pmatrix}.$$

Thus, if  $\text{rank}(\mathbf{X}) + \text{rank}(\lambda_2^{1/2}\mathbf{Q}) \geq p$  and the rows of  $\mathbf{X}$  combined with the rows of  $\lambda_2^{1/2}\mathbf{Q}$  form a linearly independent set,  $\tilde{\mathbf{C}}$  is of full rank and hence the structured elastic net is unique. In addition, this shows that even for  $p \gg n$ , in principle all features can be selected.

### 3.5 Double shrinkage

Zhou and Hastie (2005) pointed out that the elastic net effects twofold shrinkage of the coefficient vector towards zero, one produced by the  $\ell^1$ - and a second one by the  $\ell^2$ -constraint. In practice, this has the consequence that the elastic net estimates are shrunk too heavily, such that they are even outperformed by the lasso. For this reason, Zhou and Hastie (2005) suggest to undo the shrinkage induced by the  $\ell^2$ -constraint, rescaling all coefficients by the factor  $(1 + \lambda_2)$ . A similar rationale might be appropriate for our proposal, depending on the structure of  $\mathbf{\Lambda}$ . For instance, if  $\mathbf{\Lambda} = \mathbf{D}^\top \mathbf{D}$  (cf. Eq. (2.4)), then the shrinkage target of  $\mathbf{\Lambda}$  is the constant vector. In the presence of sparsity, the constant vector is roughly equal to the zero vector. Since this is already the shrinkage target of the  $\ell^1$ -constraint, we will observe double shrinkage as for the elastic net. As a remedy, one might consider to rescale the coefficients:

$$\hat{\beta}_j \leftarrow (1 + \lambda_2 l_{jj}) \hat{\beta}_j, \quad j = 1, \dots, p.$$

## 4 Consistency

The asymptotic analysis presented in this sections closely follows the ideas of Knight and Fu (2000) and Zhou (2006). Both have studied asymptotics for the lasso in linear regression under the following conditions.

(C.1) Given a sample  $S$  of size  $n$ , the data assumed to be generated according to the model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta}^* + \boldsymbol{\epsilon},$$

with  $\boldsymbol{\beta}^* = (\beta_1^*, \dots, \beta_p^*)^\top$  denoting the true parameter vector. The error terms  $\boldsymbol{\epsilon} = (\epsilon_1, \dots, \epsilon_n)^\top$  are assumed to be i.i.d. with expectation 0 and constant variance  $0 < \sigma^2 < \infty$ .

(C.2)

$$\mathbf{C}_n = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top = \frac{1}{n} \mathbf{X}_n^\top \mathbf{X}_n \rightarrow \mathbf{C} \text{ as } n \rightarrow \infty,$$

and the limit  $\mathbf{C}$  is strictly positive definite.

(C.3)

$$\frac{1}{n} \max_{1 \leq i \leq n} \mathbf{x}_i^\top \mathbf{x}_i \rightarrow 0 \text{ as } n \rightarrow \infty.$$

Conditions (C.2) and (C.3) are weak in the sense that both hold if the  $\mathbf{x}_i$  are i.i.d. with finite second moments. Conditions (C.1)-(C.3) ensure  $\sqrt{n}$ -consistency and asymptotic normality of the ordinary least squares estimator, i.e.

$$\sqrt{n}(\hat{\boldsymbol{\beta}}_n^{\text{ols}} - \boldsymbol{\beta}^*) \xrightarrow{D} N(\mathbf{0}, \sigma^2 \mathbf{C}^{-1}),$$

where here and in the following, the sub- or superscript indicates that the corresponding quantity depends on the sample size  $n$ . Using conditions (C.1)-(C.3), [Knight and Fu \(2000\)](#) proved for the lasso that  $\hat{\beta}^{\text{lasso}}$  is  $\sqrt{n}$ -consistent for  $\beta^*$  provided  $\lambda_1^n/\sqrt{n} \rightarrow \lambda_1^0 \geq 0$ . [Zhou \(2006\)](#) has shown that while taking  $\lambda_1^n = O(\sqrt{n})$  provides the optimal rate for estimation, it leads to inconsistent feature selection. Define the active set as  $A = \{j : \beta_j^* \neq 0\}$  and  $A^c = \{1, \dots, p\} \setminus A$  and let  $\delta$  be an estimation procedure producing an estimate  $\hat{\beta}^\delta$ . Then  $\delta$  is said to consistent in feature selection if

$$\begin{aligned} \lim_{n \rightarrow \infty} P(\hat{\beta}_{j,n}^\delta \neq 0) &= 1 \quad \text{for } j \in A, \\ \lim_{n \rightarrow \infty} P(\hat{\beta}_{j,n}^\delta = 0) &= 1 \quad \text{for } j \in A^c. \end{aligned}$$

Moreover, [Zhou \(2006\)](#) and [Zhao and Yu \(2006\)](#) have shown that if  $\lambda_1^n/n \rightarrow 0$  and  $\lambda_1^n/\sqrt{n} \rightarrow \infty$ , the lasso has to satisfy a nontrivial condition, the so-called 'irrepresentable condition', to be selection consistent. [Zhou \(2006\)](#) proposed the adaptive lasso, a two-step estimation procedure, to fix this deficiency. In the following, these results will be adapted to the presence of a second quadratic penalty term.

**Theorem 1.** *Let conditions (C.1)-(C.3) hold. Define*

$$\hat{\beta}_n = \underset{\beta}{\operatorname{argmin}} \|\mathbf{y}_n - \mathbf{X}_n \beta\|^2 + \lambda_1^n \|\beta\|_1 + \lambda_2^n \beta^\top \mathbf{\Lambda} \beta.$$

*Assume that  $\lambda_1^n/\sqrt{n} \rightarrow \lambda_1^0 \geq 0$  and  $\lambda_2^n/\sqrt{n} \rightarrow \lambda_2^0 \geq 0$ . Consider the random function*

$$\begin{aligned} V(\mathbf{u}) &= -2\mathbf{u}^\top \mathbf{w} + \mathbf{u}^\top \mathbf{C} \mathbf{u} \\ &\quad + \lambda_1^0 \sum_{j=1}^p u_j \operatorname{sign}(\beta_j^*) I(\beta_j^* \neq 0) + |u_j| I(\beta_j^* = 0) \\ &\quad + 2\lambda_2^0 \mathbf{u}^\top \mathbf{\Lambda} \beta^*, \quad \mathbf{w} \sim N(\mathbf{0}, \sigma^2 \mathbf{C}). \end{aligned}$$

*Then  $\sqrt{n}(\hat{\beta}_n - \beta^*) \xrightarrow{D} \operatorname{argmin} V(\mathbf{u})$ .*

Theorem 1 is analogous to Theorem 2 in [Knight and Fu \(2000\)](#) and establishes  $\sqrt{n}$ -consistency of  $\hat{\beta}_n$ , provided  $\lambda_1^n$  and  $\lambda_2^n$  are  $O(\sqrt{n})$ . Theorem 1 admits a straightforward extension to the class of generalized linear models (cf. Eq. (3.2)). Let the true model be defined by

$$E[Y|\mathbb{X} = \mathbf{x}] = b'(f(\mathbf{x}; \beta^*)), \quad f(\mathbf{x}) = \mathbf{x}^\top \beta^*.$$

For the sake of a clearer presentation, we assume that  $\beta_0^* = 0$ . We study the estimator

$$\hat{\beta}_n = \underset{\beta}{\operatorname{argmin}} 2\phi^{-1} \sum_{i=1}^n b(f(\mathbf{x}_i; \beta)) - y_i f(\mathbf{x}_i; \beta) + \lambda_1^n \|\beta\|_1 + \lambda_2^n \beta^\top \mathbf{\Lambda} \beta. \quad (4.1)$$

We work with the following regularity conditions.

(G.1) The expected Fisher information

$$\mathcal{I} = E[\phi^{-1} b''(f(\mathbb{X}); \beta^*) \mathbb{X} \mathbb{X}^\top]$$

is finite and strictly positive definite.

(G.2) There exists a function  $M$  and an open neighbourhood  $U$  of  $\beta^*$  such that for all  $\beta \in U$

$$|b'''(f(\mathbf{x}); \beta)| \leq M(\mathbf{x}) < \infty \text{ for all } \mathbf{x},$$

and

$$\mathbb{E}[M(\mathbb{X}) |X_j X_k X_l|] < \infty \quad \forall 1 \leq j, k, l \leq p.$$

Condition (G.2) is necessary for a Taylor expansion argument, as explained in the Appendix.

**Theorem 2.** *Let conditions (G.1) and (G.2) hold, consider the estimator (4.1) and let  $\lambda_1^n/\sqrt{n} \rightarrow \lambda_1^0 \geq 0$  and  $\lambda_2^n/\sqrt{n} \rightarrow \lambda_2^0 \geq 0$ . Consider the random function*

$$\begin{aligned} W(\mathbf{u}) = & -2\mathbf{u}^\top \mathbf{w} + \mathbf{u}^\top \mathcal{I} \mathbf{u} \\ & + \lambda_1^0 \sum_{j=1}^p u_j \text{sign}(\beta_j^*) I(\beta_j^* \neq 0) + |u_j| I(\beta_j^* = 0) \\ & + 2\lambda_2^0 \mathbf{u}^\top \mathbf{\Lambda} \beta^*, \quad \mathbf{w} \sim N(\mathbf{0}, \mathcal{I}). \end{aligned}$$

Then  $\sqrt{n}(\hat{\beta}_n - \beta^*) \xrightarrow{D} \text{argmin } W(\mathbf{u})$ .

Now let us turn to the question of selection consistency. In the situation of Theorem 1, if  $\lambda_1^n$  and  $\lambda_2^n$  both are  $O(\sqrt{n})$ , then  $\hat{\beta}_n$  cannot be selection consistent. To see this, let  $\hat{\mathbf{u}}$  denote the minimizer of  $V(\mathbf{u})$ . Without loss of generality, let us assume that  $\beta^* = ([\beta_A^*]^\top, [\beta_{A^c}^*]^\top)^\top$ , where here and in the following, the subscripts  $A$  and  $A^c$  refer to active and inactive set, respectively. Selection consistency requires that  $\hat{\mathbf{u}}$  is of the form  $\hat{\mathbf{u}} = (\hat{\mathbf{u}}_A^\top, \mathbf{0}^\top)^\top$ . Evaluation of  $V(\mathbf{u})$  at  $\hat{\mathbf{u}}$  yields:

$$V(\hat{\mathbf{u}}) = -2\hat{\mathbf{u}}_A^\top \mathbf{w}_A + \hat{\mathbf{u}}_A^\top \mathbf{C}_A \hat{\mathbf{u}}_A + \lambda_1^0 \hat{\mathbf{u}}_A^\top \mathbf{s}_A + 2\lambda_2^0 \hat{\mathbf{u}}_A^\top \mathbf{\Lambda}_A \beta_A^*, \quad \mathbf{s}_A = (\text{sign}(\beta_j^*), j \in A)^\top.$$

Since  $\hat{\mathbf{u}}_A$  is a minimizer, differentiation yields

$$\hat{\mathbf{u}}_A = \mathbf{C}_A^{-1} \left( \mathbf{w}_A - \frac{\lambda_1^0}{2} \mathbf{s}_A - \lambda_2^0 \mathbf{\Lambda}_A \beta_A^* \right). \quad (4.2)$$

Now consider the partitioning schemes

$$\mathbf{C} = \begin{pmatrix} \mathbf{C}_A & \mathbf{C}_{AA^c} \\ \mathbf{C}_{A^c A} & \mathbf{C}_{A^c A^c} \end{pmatrix} \quad \text{and} \quad \mathbf{\Lambda} = \begin{pmatrix} \mathbf{\Lambda}_A & \mathbf{\Lambda}_{AA^c} \\ \mathbf{\Lambda}_{A^c A} & \mathbf{\Lambda}_{A^c A^c} \end{pmatrix}. \quad (4.3)$$

Since  $\hat{\mathbf{u}}_{A^c} = \mathbf{0}$ , the Karush-Kuhn-Tucker (KKT) conditions imply that

$$|\lambda_2^0 \mathbf{\Lambda}_{A^c A} \beta_A^* + \mathbf{C}_{A^c A} \hat{\mathbf{u}}_A - \mathbf{w}_{A^c}| \leq \frac{\lambda_1^0}{2} \mathbf{1}, \quad (4.4)$$

where the inequality is interpreted componentwise. Substituting the right hand side of Eq. (4.2) into Eq. (4.4) and rearranging terms, one obtains

$$\left| [\mathbf{\Lambda}_{A^c A} - \mathbf{C}_{A^c A} \mathbf{C}_A^{-1} \mathbf{\Lambda}_A] \lambda_2^0 \beta_A^* + \mathbf{C}_{A^c A} \mathbf{C}_A^{-1} \left( \mathbf{w}_A - \frac{\lambda_1^0}{2} \mathbf{s}_A \right) - \mathbf{w}_{A^c} \right| \leq \frac{\lambda_1^0}{2} \mathbf{1}. \quad (4.5)$$

Since  $\mathbf{w}$  is random quantity, the system of inequalities only holds with a certain probability, implying inconsistency of selection in general. In addition, the system (4.5) shows that selection consistency depends on  $\mathbf{C}$ ,  $\mathbf{\Lambda}$  and also  $\beta_A^*$ . Selection consistency can be achieved if one lets  $\lambda_1^n, \lambda_2^n$  grow more strongly and if the quantities  $\mathbf{C}$ ,  $\mathbf{\Lambda}$  and  $\beta_A^*$  fulfill a nontrivial condition, which can be seen as analog to the irrepresentable condition of the lasso.

**Theorem 3.** *In the setup of Theorem 1, let  $\lambda_1^n/n \rightarrow 0$ ,  $\lambda_1^n/\sqrt{n} \rightarrow \infty$ ,  $\lambda_2^n/\lambda_1^n \rightarrow R$ ,  $0 < R < \infty$ . Then, if selection consistency holds, the following condition must be fulfilled: there exists a sign vector  $\mathbf{s}_A$  such that*

$$| -\mathbf{C}_{A^c A} \mathbf{C}_A^{-1} (\mathbf{s}_A + 2R\mathbf{\Lambda}_A \boldsymbol{\beta}_A^*) + 2R\mathbf{\Lambda}_{A^c A} \boldsymbol{\beta}_A^* | \leq \mathbf{1},$$

where the inequality is interpreted componentwise.

While this condition is interesting from a theoretical point of view, it is impossible to check in practice, since  $\boldsymbol{\beta}_A^*$  is unknown.

Selection consistency can be achieved by a two-step estimation strategy introduced in Zhou (2006) under the name adaptive lasso, which replaces  $\ell^1$ -regularization uniform in  $\beta_j$ ,  $j = 1, \dots, p$ , by a weighted variant  $J(\boldsymbol{\beta}) = \sum_{j=1}^p \omega_j |\beta_j|$ , where the weights  $\{\omega_j\}_{j=1}^p$  are determined adaptively as function of an 'initial estimator'  $\hat{\boldsymbol{\beta}}^{\text{init}}$ :

$$\omega_j = |\hat{\beta}_j^{\text{init}}|^{-\gamma}, \quad \gamma > 0, \quad j = 1, \dots, p. \quad (4.6)$$

In terms of selection consistency this strategy turns out to be favourable for our proposal, too.

**Theorem 4.** *In the setup of Theorem 1, define*

$$\hat{\boldsymbol{\beta}}_n^{\text{adaptive}} = \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \|\mathbf{y}_n - \mathbf{X}_n \boldsymbol{\beta}\|^2 + \lambda_1^n \sum_{j=1}^p \omega_j |\beta_j| + \lambda_2^n \boldsymbol{\beta}^\top \mathbf{\Lambda} \boldsymbol{\beta},$$

where the weights are as in Eq. (4.6), and suppose that the initial estimator satisfies

$$r_n(\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}^*) = O_P(1), \quad r_n \rightarrow \infty \text{ as } n \rightarrow \infty.$$

Furthermore, suppose that

$$r_n^\gamma \lambda_1^n n^{-1/2} \rightarrow \infty, \quad \lambda_1^n n^{-1/2} \rightarrow 0, \quad \lambda_2^n n^{-1/2} \rightarrow \lambda_2^0 \geq 0$$

as  $n \rightarrow \infty$ . Then

- (1)  $\sqrt{n}(\hat{\boldsymbol{\beta}}_{A,n}^{\text{adaptive}} - \boldsymbol{\beta}_A^*) \xrightarrow{D} N(-\lambda_2^0 \mathbf{C}_A^{-1} \mathbf{\Lambda}_A \boldsymbol{\beta}_A^*, \mathbf{C}_A^{-1}),$
- (2)  $\lim_{n \rightarrow \infty} P(\hat{\boldsymbol{\beta}}_{A^c,n}^{\text{adaptive}} = \mathbf{0}) = 1.$

Theorem 4 implies that the adaptive structured elastic net  $\hat{\boldsymbol{\beta}}^{\text{adaptive}}$  is an oracle estimation procedure (Fan and Li (2001)) if the bias term in (1) vanishes, which is the case if  $\boldsymbol{\beta}_A^*$  resides in the nullspace of  $\mathbf{\Lambda}_A$ . Interestingly, if  $\mathbf{\Lambda}$  equals the combinatorial graph Laplacian (cf. Section 2.1), this happens if and only if  $\boldsymbol{\beta}_A^*$  has constant entries and  $A$  specifies a connected component in the underlying graph.

Concerning the choice of the initial estimator, the ridge estimator has worked well for us in practice, provided the ridge parameter is chosen appropriately. While  $\gamma$  may be treated as a tuning parameter, we have set  $\gamma$  equal to 1 in all our data analyses. Finally, we remark that while Theorem 4 applies to linear regression, it can be extended to hold for generalized linear models, in a similar way as the extension of Theorem 1 to Theorem 2.

## 5 Computation

This section discusses aspects concerning computation and model selection for the structured elastic net estimator when the loss function is the negative log-likelihood of a generalized linear model (3.1).

## 5.1 Data augmentation

From the discussions in Subsections 3.3 and 3.4, it follows that the structured elastic net for squared loss, assuming centered data, can be recast as lasso on augmented data

$$\widetilde{\mathbf{X}} = \begin{pmatrix} \mathbf{X} \\ \lambda_2^{1/2} \mathbf{Q} \end{pmatrix}_{(n+p) \times p}, \quad \widetilde{\mathbf{y}} = \begin{pmatrix} \mathbf{y} \\ \mathbf{0} \end{pmatrix}_{(n+p) \times 1}, \quad \mathbf{\Lambda} = \mathbf{Q}^\top \mathbf{Q},$$

and hence algorithms available for computing the lasso, notably LARS (Efron et al. (2004)) may be applied, which computes for fixed  $\lambda_2$  and varying  $\lambda_1$  the piecewise linear solution path  $\widehat{\boldsymbol{\beta}}(\lambda_1; \lambda_2)$ . In order to fit arbitrary regularized generalized linear models, the augmented data representation has to be modified. Without regularization, estimators in generalized linear models are obtained by iteratively computing weighted least squares estimators:

$$\begin{aligned} \begin{pmatrix} \widehat{\beta}_0^{(k+1)} \\ \widehat{\boldsymbol{\beta}}^{(k+1)} \end{pmatrix} &= \left( [\mathbf{1} \ \mathbf{X}]^\top \mathbf{W}^{(k)} [\mathbf{1} \ \mathbf{X}] \right)^{-1} [\mathbf{1} \ \mathbf{X}]^\top \mathbf{W}^{(k)} \mathbf{z}^{(k)}, \\ \mathbf{z}^{(k)} &= \mathbf{f}^{(k)} + [\mathbf{W}^{(k)}]^{-1} (\mathbf{y} - \boldsymbol{\mu}^{(k)}), \\ \mathbf{f}^{(k)} &= (f_1^{(k)}, \dots, f_n^{(k)})^\top, \quad f_i^{(k)} = \widehat{\beta}_0^{(k)} + \mathbf{x}_i^\top \widehat{\boldsymbol{\beta}}^{(k)}, \quad i = 1, \dots, n, \\ \boldsymbol{\mu}^{(k)} &= (\mu_1^{(k)}, \dots, \mu_n^{(k)})^\top, \quad \mu_i^{(k)} = b'(f_i^{(k)}), \quad i = 1, \dots, n, \\ \mathbf{W}^{(k)} &= \text{diag}(w_1^{(k)}, \dots, w_n^{(k)}), \quad w_i^{(k)} = \phi^{-1} b''(f_i^{(k)}), \quad i = 1, \dots, n. \end{aligned} \tag{5.1}$$

Note that the design matrix additionally includes a constant term  $\mathbf{1}$ . Turning back to the structured elastic net, an adaptation of the augmented data approach iteratively determines

$$\begin{pmatrix} \widehat{\beta}_0^{(k+1)} \\ \widehat{\boldsymbol{\beta}}^{(k+1)} \end{pmatrix} = \underset{(\beta_0, \boldsymbol{\beta})}{\operatorname{argmin}} \sum_{i=1}^{n+p} \widetilde{w}_i^{(k)} \left( \widetilde{z}_i^{(k)} - \widetilde{\mathbf{x}}_i^\top \begin{pmatrix} \beta_0 \\ \boldsymbol{\beta} \end{pmatrix} \right)^2 + \lambda_1 \|\boldsymbol{\beta}\|_1,$$

with

$$\begin{aligned} \widetilde{w}_i^{(k)} &= w_i^{(k)}, \quad i = 1, \dots, n, \text{ as in Eq. (5.1),} \quad \widetilde{w}_i^{(k)} = 1, \quad i = (n+1), \dots, (n+p), \\ \widetilde{z}_i^{(k)} &= z_i^{(k)}, \quad i = 1, \dots, n, \text{ as in Eq. (5.1),} \quad \widetilde{z}_i^{(k)} = 0, \quad i = (n+1), \dots, (n+p), \\ \widetilde{\mathbf{x}}_i &= (\mathbf{1} \ \mathbf{x}_i^\top)^\top, \quad i = 1, \dots, n, \quad \widetilde{\mathbf{x}}_i = (0 \ \sqrt{\lambda_2} \mathbf{q}_i^\top)^\top, \quad i = (n+1), \dots, (n+p), \end{aligned}$$

with  $\mathbf{q}_i^\top$  denoting the  $i$ -th row of  $\mathbf{Q}$ .

## 5.2 Cyclical coordinate descent

Recently, cyclical coordinate descent (CCD) approaches for obtaining the lasso solution in generalized linear models have gained much popularity, e.g. Genkin et al. (2007), Wu and Lange (2008), Friedman et al. (2007). The latter have studied convex functions of the form

$$F(\boldsymbol{\beta}) = S(\boldsymbol{\beta}) + J(\boldsymbol{\beta}), \quad \boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^\top, \quad J(\boldsymbol{\beta}) = \sum_{j=1}^p \iota_j(\beta_j), \tag{5.2}$$

to be minimized w.r.t.  $\boldsymbol{\beta}$ . We suppose that  $F(\boldsymbol{\beta})$  is convex, that  $S(\boldsymbol{\beta})$  is smooth and that  $J(\boldsymbol{\beta})$  is continuous and separable in the  $\beta_j$ . CCD optimizes one coordinate at a time, with

the remaining ones kept fixed. Convergence analysis for this setup is studied in [Tseng \(2001\)](#). As reported in [Friedman et al. \(2007\)](#), the CCD approach turns out to be remarkably efficient as well as versatile, allowing the computation of the lasso, the elastic net and several related procedures along a fine grid of  $\lambda_1$ -values, using the current estimate as 'warm start' for the next grid point. While CCD could be applied to the augmented data representation of the previous subsection to compute the structured elastic net estimator, we prefer a CCD algorithm directly adapted to our specific problem, without the need to determine the root  $\mathbf{Q}$ . Note that such an algorithm is possible, since the structured elastic net criterion matches the structure of  $F(\boldsymbol{\beta})$  in Eq. (5.2) with  $\iota_j = \lambda_1 |\beta_j|$ ,  $j = 1, \dots, p$ . We first state the algorithm for squared loss, and straightforward modifications admit the extension to generalized linear models.

Consider the structured elastic net estimator for  $L(y, f(\mathbf{x}; \beta_0, \boldsymbol{\beta})) = (y - \beta_0 - \mathbf{x}^\top \boldsymbol{\beta})^2$ . The aim is to determine an estimate  $\hat{\beta}_j$  given  $\hat{\beta}_0, \hat{\boldsymbol{\beta}}_{-j}$ ,  $\hat{\boldsymbol{\beta}}_{-j} = (\hat{\beta}_1, \dots, \hat{\beta}_{j-1}, \hat{\beta}_{j+1}, \dots, \hat{\beta}_p)^\top$ . Likewise, the  $j$ -th row of  $\mathbf{A}$  can be divided into  $l_{jj}$  and  $\mathbf{l}_{-j} = (l_{jj'})_{j' \neq j}^\top$ . The KKT conditions imply that if  $\hat{\beta}_j = 0$

$$\left| -2\mathbf{X}_j^\top \left( \mathbf{y} - \hat{\beta}_0 - \sum_{j' \neq j} \mathbf{X}_{j'} \hat{\beta}_{j'} \right) + 2\lambda_2 \mathbf{l}_{-j}^\top \hat{\boldsymbol{\beta}}_{-j} \right| \leq \lambda_1, \quad (5.3)$$

and otherwise,

$$-2\mathbf{X}_j^\top (\mathbf{y} - \hat{\beta}_0 - \mathbf{X} \hat{\boldsymbol{\beta}}) + 2\lambda_2 l_{jj} \hat{\beta}_j + 2\lambda_2 \mathbf{l}_{-j}^\top \hat{\boldsymbol{\beta}}_{-j} + \lambda_1 \text{sign}(\hat{\beta}_j) = 0. \quad (5.4)$$

Solving for  $\hat{\beta}_j$  in Eq. (5.4) yields

$$\begin{aligned} \hat{\beta}_j &= \frac{\mathbf{X}_j^\top (\mathbf{r}_j - \mathbf{l}_{-j}^\top \hat{\boldsymbol{\beta}}_{-j}) - \frac{\lambda_1}{2} \text{sign}(\hat{\beta}_j)}{\|\mathbf{X}_j\|^2 + \lambda_2 l_{jj}}, \\ \mathbf{r}_j &= \mathbf{y} - \hat{\beta}_0 - \sum_{j' \neq j} \mathbf{X}_{j'} \hat{\beta}_{j'}. \end{aligned} \quad (5.5)$$

Combining Eq. (5.3) with Eq. (5.4) yields the update formula

$$\hat{\beta}_j \leftarrow \frac{\left[ \left| \mathbf{X}_j^\top \mathbf{r}_j - \mathbf{l}_{-j}^\top \hat{\boldsymbol{\beta}}_{-j} \right| - \frac{\lambda_1}{2} \right]_+ \text{sign}(\mathbf{X}_j^\top \mathbf{r}_j - \mathbf{l}_{-j}^\top \hat{\boldsymbol{\beta}}_{-j})}{\|\mathbf{X}_j\|^2 + \lambda_2 l_{jj}}, \quad j = 1, \dots, p, \quad [z]_+ = \max(z, 0).$$

For the intercept, we have the update

$$\hat{\beta}_0 \leftarrow \frac{\sum_{i=1}^n y_i - \sum_{j=1}^p \mathbf{X}_j \hat{\beta}_j}{n}.$$

For generalized linear models, the algorithm involves one outer loop in which the quantities of iteratively weighted least squares (5.1) are updated, and one inner loop for CCD as just described. To be more precise, given estimates  $\hat{\beta}_0^{(k)}, \hat{\boldsymbol{\beta}}^{(k)}$ , we compute the  $\mathbf{W}^{(k)}$  and  $\mathbf{z}^{(k)}$  as given in Eq. (5.1). Defining the working response  $\tilde{\mathbf{z}}^{(k)} = [\mathbf{W}^{(k)}]^{1/2} \mathbf{z}^{(k)}$ ,  $\tilde{\mathbf{X}}^{(k)} = [\mathbf{W}^{(k)}]^{1/2} \mathbf{X}$  and a modified intercept  $\tilde{\mathbf{1}}^{(k)} = [\tilde{\mathbf{W}}^{(k)}]^{1/2} \mathbf{1}$ , the data  $(\tilde{\mathbf{z}}^{(k)}, [\tilde{\mathbf{1}}^{(k)} \tilde{\mathbf{X}}^{(k)}])$  can be plugged into the CCD algorithm for squared loss. Once an inner loop has converged, the new estimates  $\hat{\beta}_0^{(k+1)}, \hat{\boldsymbol{\beta}}^{(k+1)}$  are used to obtain  $(\tilde{\mathbf{z}}^{(k+1)}, [\tilde{\mathbf{1}}^{(k+1)} \tilde{\mathbf{X}}^{(k+1)}])$ .



### 5.3 Goeman's algorithm

The third algorithm we provide is an adaptation of a recent proposal in [Goeman \(2007\)](#), which combines full gradient descent- with Newton-Raphson-steps in subdomains where the gradient of the structured elastic net objective function is continuous. For a loss function of the type (3.1), the gradient and the Hessian of the differentiable part of the structured elastic net criterion are given by

$$\begin{aligned} \begin{pmatrix} \nabla_{\beta_0} \\ \nabla_{\beta} \end{pmatrix} &= \begin{pmatrix} -\mathbf{1}^\top(\mathbf{y} - \boldsymbol{\mu}(\beta_0, \boldsymbol{\beta})) \\ -\mathbf{X}^\top(\mathbf{y} - \boldsymbol{\mu}(\beta_0, \boldsymbol{\beta})) + 2\lambda_2\boldsymbol{\Lambda} \end{pmatrix}, \\ \nabla_{\beta_0, \boldsymbol{\beta}}^2 &= \begin{pmatrix} \nabla_{\beta_0}^2 & \frac{\partial}{\partial \beta^\top} \nabla_{\beta_0} \\ \frac{\partial}{\partial \beta_0} \nabla_{\beta} & \nabla_{\beta}^2 \end{pmatrix} = \begin{pmatrix} \mathbf{1}^\top \mathbf{W}(\beta_0, \boldsymbol{\beta}) \mathbf{1} & \mathbf{1}^\top \mathbf{W}(\beta_0, \boldsymbol{\beta}) \mathbf{X} \\ \mathbf{X}^\top \mathbf{W}(\beta_0, \boldsymbol{\beta}) \mathbf{1} & \mathbf{X}^\top \mathbf{W}(\beta_0, \boldsymbol{\beta}) \mathbf{X} + 2\lambda_2\boldsymbol{\Lambda} \end{pmatrix}. \end{aligned} \quad (5.6)$$

Let the superscript  $(k)$  refer to the iteration counter and let  $A$  denote the currently active set, i.e.  $A = \{j \in \{1, \dots, p\} : \hat{\beta}_j^{(k)} \neq 0\}$ . In principle, the idea is to take gradient descent steps until it is favourable to switch to Newton-Raphson. For the former, one looks for a suitable descent direction  $\mathbf{v}$ . For  $j \in A$ , the  $\ell^1$ -part of the regularizer is differentiable as well: combined with the gradient in Eq. (5.6), this yields descent directions

$$v_0 \leftarrow -\nabla_{\hat{\beta}_0^{(k)}}, \quad \mathbf{v}_A \leftarrow -[\nabla_{\hat{\boldsymbol{\beta}}^{(k)}}]_A - \lambda_1 \text{sign}(\hat{\boldsymbol{\beta}}_A),$$

where the evaluation of the sign function here and in the following is understood componentwise. For  $j \in A^c$ , we have to distinguish two cases. If  $||\nabla_{\hat{\boldsymbol{\beta}}^{(k)}}]_j| > \lambda_1$ , one enlarges the active set, i.e.  $A \leftarrow A \cup \{j\}$ , and sets  $v_j \leftarrow -[\nabla_{\hat{\boldsymbol{\beta}}^{(k)}}]_j - \lambda_1 \text{sign}([\nabla_{\hat{\boldsymbol{\beta}}^{(k)}}]_j)$ . Otherwise, there is no move for coordinate  $j$ :  $v_j \leftarrow 0$ . We define  $\mathbf{v} = (v_1, \dots, v_p)^\top$ . The next issue to consider is the determination of an appropriate step length  $t$  for the gradient descent update  $\hat{\boldsymbol{\beta}}^{(k+1)} \leftarrow \hat{\boldsymbol{\beta}}^{(k)} + t\mathbf{v}$ . One has to take care that the stepsize is sufficiently small in order to remain within a subdomain of gradient continuity, i.e. one has to ensure that  $\text{sign}(\hat{\boldsymbol{\beta}}^{(k+1)}) = \text{sign}(\hat{\boldsymbol{\beta}}^{(k)})$ . This naturally imposes an upper bound on the step size, denoted by  $t^{\text{edge}}$ :

$$0 < t^{\text{edge}} = \min_{j \in A} \left\{ -\frac{\hat{\beta}_j^{(k)}}{v_j} : \text{sign}(\hat{\beta}_j^{(k)}) = -\text{sign}(v_j) \right\}. \quad (5.7)$$

In a subdomain of gradient continuity, assuming  $t < t^{\text{edge}}$ , one can perform a quadratic Taylor approximation of the objective function  $F(\beta_0, \boldsymbol{\beta})$  around the current estimates  $\hat{\beta}_0^{(k)}, \hat{\boldsymbol{\beta}}^{(k)}$ :

$$F(\hat{\beta}_0^{(k)}, \hat{\boldsymbol{\beta}}^{(k)} + t\mathbf{v}) \approx F(\hat{\beta}_0^{(k)}, \hat{\boldsymbol{\beta}}^{(k)}) - t\mathbf{v}^\top \mathbf{v} + \frac{1}{2}t^2 \mathbf{v}^\top \nabla_{\hat{\boldsymbol{\beta}}^{(k)}}^2 \mathbf{v}. \quad (5.8)$$

Differentiating the r.h.s. of approximation (5.8) with respect to  $t$  and setting the result equal to zero, it follows that the minimum of the Taylor approximation is achieved by choosing the stepsize

$$t^{\text{opt}} = \frac{\mathbf{v}^\top \mathbf{v}}{\mathbf{v}^\top \nabla_{\hat{\boldsymbol{\beta}}^{(k)}}^2 \mathbf{v}},$$

provided the denominator is nonzero. At this place, one checks whether  $t^{\text{opt}} < t^{\text{edge}}$ . If this is the case, the gradient descent update is performed with step size  $t = t^{\text{opt}}$ , i.e.

$$\hat{\boldsymbol{\beta}}^{(k+1)} \leftarrow \hat{\boldsymbol{\beta}}^{(k)} + t^{\text{opt}} \mathbf{v}.$$

Otherwise, the active set is reduced:  $A \leftarrow A \setminus \{j_0\}$  and  $\hat{\beta}_{j_0}^{(k+1)} \leftarrow 0$ , where  $j_0$  is the index for which the minimum in Equation (5.7) is attained, and the gradient descent step is performed with stepsize  $t^{\text{edge}}$  for the reduced active set. For the intercept, one may use  $\hat{\beta}_0^{(k+1)} = \hat{\beta}_0^{(k)} + t_0 v_0$ ,  $t_0 = 1/\nabla_{\hat{\beta}_0}^2$ .

Near the minimum, Newton-Raphson can considerably speed up convergence. Newton-Raphson-steps are only possible within a single subdomain of gradient continuity. As Newton-Raphson-steps are computationally expensive, they should be avoided if they are likely to fail. The Newton-Raphson-update restricted to the active set is given by

$$\hat{\beta}_A^{(k+1)} \leftarrow \hat{\beta}_A^{(k)} - [\nabla_{\hat{\beta}^{(k)}}^2]_A^{-1} \left\{ [\nabla_{\hat{\beta}^{(k)}}]_A + \lambda_1 \text{sign}(\hat{\beta}_A^{(k)}) \right\}. \quad (5.9)$$

The Newton-Raphson-step is accepted only if  $\text{sign}(\hat{\beta}_A^{(k+1)}) = \text{sign}(\hat{\beta}_A^{(k)})$ . Otherwise, a subdomain of gradient continuity has been left. Following Goeman (2007), a Newton-Raphson step should be attempted only if  $t^{\text{opt}} < t^{\text{edge}}$  in a preceding gradient descent step. Moreover, Newton-Raphson-steps are neither favourable in the very first iterations nor when the active set has not yet changed since the last failed attempt.

If  $p \gg n$ , it is possible that the size of the active set exceeds by far the sample size. In this case, the Newton-Raphson-steps given by the update (5.9) are inefficient. The matrix inversion requires  $O(|A|^3)$  operations, but it is immediately seen that a reduction to  $O(n^3)$  operations can be achieved by reparametrization. We set  $\hat{\beta}_A^{(k)} = \mathbf{X}_A^\top \hat{\gamma}^{(k)}$ , so that the gradient has the form

$$\nabla_{\hat{\gamma}^{(k)}} = -\mathbf{X}_A \mathbf{X}_A^\top (\mathbf{y} - \boldsymbol{\mu}) + \lambda_1 \mathbf{X}_A \text{sign}(\mathbf{X}_A^\top \hat{\gamma}^{(k)}) + 2\lambda_2 \mathbf{X}_A \boldsymbol{\Lambda}_A \mathbf{X}_A^\top \hat{\gamma}^{(k)} = \mathbf{X}_A [\nabla_{\hat{\beta}^{(k)}}]_A,$$

and for the Hessian, one obtains

$$\nabla_{\hat{\gamma}^{(k)}}^2 = \mathbf{X}_A \mathbf{X}_A^\top \mathbf{W}_A \mathbf{X}_A \mathbf{X}_A^\top + 2\lambda_2 \mathbf{X}_A \boldsymbol{\Lambda}_A \mathbf{X}_A^\top = \mathbf{X}_A [\nabla_{\hat{\beta}^{(k)}}^2]_A \mathbf{X}_A^\top.$$

For the Newton-Raphson-steps, this has the following algorithmic implications.

- One checks if  $|A| \gg n$ , where ' $\gg$ ' means that the ratio  $|A|/n$  exceeds a factor related to the specific problem.
- If yes, one reparametrizes  $\hat{\beta}_A^{(k)}$  to  $\hat{\gamma}^{(k)}$  by solving the linear equation  $\mathbf{X}_A \mathbf{X}_A^\top \hat{\gamma}^{(k)} = \mathbf{X}_A \hat{\beta}_A^{(k)}$ . Otherwise, one proceeds according to the update (5.9).
- One applies Newton-Raphson to  $\hat{\gamma}^{(k)}$ , using  $\nabla_{\hat{\gamma}^{(k)}}$ ,  $\nabla_{\hat{\gamma}^{(k)}}^2$  given above to obtain  $\hat{\gamma}^{(k+1)}$ .
- One backtransforms  $\hat{\beta}_A^{(k+1)} = \mathbf{X}_A^\top \hat{\gamma}^{(k+1)}$ .
- The Newton-Raphson step is accepted only if  $\text{sign}(\hat{\beta}_A^{(k+1)}) = \text{sign}(\hat{\beta}_A^{(k)})$ .

## 5.4 Comparison

Cyclical coordinate descent is conceptually simple and involves only vector operations, also avoiding the direct storage of the matrix  $\mathbf{X}^\top \mathbf{X}$ . If organized properly, one inner loop involves  $O(p(n+p))$  operations. Although we do not give a proof of convergence, the algorithm is well-founded since it can be embedded into the framework of Tseng (2001). This is contrary to Goeman's algorithm, for which no convergence analysis has been performed yet. Several parts resort to heuristics, but it generally works well in practice and has turned out to be faster than CCD, which can be rather slow for large  $p$ .

## 5.5 Degrees of freedom

The 'degrees of freedom' of some estimation procedure  $\delta$  are an integral part of model selection criteria such as GCV,  $C_p$  or the AIC. Therefore, it would be desirable to determine the degrees of freedom of the structured elastic net as a function of the hyperparameters  $\lambda_1$  and  $\lambda_2$ . Given a sample  $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$  such that  $\boldsymbol{\mu} = \mathbb{E}[\mathbf{y}|\mathbf{X}]$  and  $\text{var}[\mathbf{y}|\mathbf{X}] = \sigma^2 \mathbf{I}$ , denote by  $\hat{\boldsymbol{\mu}}$  the fit when applying  $\delta$  to  $S$ . In the framework of Stein's unbiased risk estimation (Stein (1981)), the degrees of freedom of  $\hat{\boldsymbol{\mu}}$  are given by

$$\text{df}(\hat{\boldsymbol{\mu}}) = \sum_{i=1}^n \text{cov}(\hat{\mu}_i, y_i) / \sigma^2. \quad (5.10)$$

This definition can easily be evaluated if  $\hat{\boldsymbol{\mu}} = \mathbf{A}\mathbf{y}$  for some matrix  $\mathbf{A}$  independent of  $\mathbf{y}$ , in which case  $\text{df}(\hat{\boldsymbol{\mu}}) = \text{tr}(\mathbf{A})$ . The lasso fit is a nonlinear function of  $\mathbf{y}$ , which considerably complicates the evaluation of (5.10). A heuristic due to Tibshirani (1996), see also Fan and Li (2001), represents the lasso fit for linear models as weighted ridge fit:

$$\hat{\boldsymbol{\mu}} \approx [\mathbf{1} \ \mathbf{X}]([\mathbf{1} \ \mathbf{X}]^\top [\mathbf{1} \ \mathbf{X}] + \lambda_1 \boldsymbol{\Omega})^{-1} [\mathbf{1} \ \mathbf{X}]^\top \mathbf{y} = \mathbf{A}\mathbf{y},$$

where  $\boldsymbol{\Omega} = \text{diag}(0, (\hat{\beta}_1^{\text{lasso}})^2 / |\hat{\beta}_1^{\text{lasso}}|, \dots, (\hat{\beta}_p^{\text{lasso}})^2 / |\hat{\beta}_p^{\text{lasso}}|)$ , with the convention  $0^{-1} = 0$ . The degrees of freedom according to Eq. (5.10) are then computed as the trace of  $\mathbf{A}$ . This heuristic can be modified for the structured elastic net fit via

$$\hat{\boldsymbol{\mu}} = [\mathbf{1} \ \mathbf{X}]([\mathbf{1} \ \mathbf{X}]^\top [\mathbf{1} \ \mathbf{X}] + \lambda_1 \boldsymbol{\Omega} + \lambda_2 \tilde{\boldsymbol{\Lambda}})^{-1} [\mathbf{1} \ \mathbf{X}]^\top \mathbf{y}, \quad \tilde{\boldsymbol{\Lambda}} = \begin{pmatrix} 0 & \mathbf{0}^\top \\ \mathbf{0} & \boldsymbol{\Lambda} \end{pmatrix}.$$

A similar formula can be applied for the adaptive structured elastic net by rescaling the entries of  $\boldsymbol{\Omega}$ . For generalized linear models, we can make use of the iteratively weighted least squares approximation (5.1):

$$\hat{\boldsymbol{\mu}} \approx [\mathbf{W}(\hat{\beta}_0, \hat{\boldsymbol{\beta}})]^{1/2} [\mathbf{1} \ \mathbf{X}]([\mathbf{1} \ \mathbf{X}]^\top \mathbf{W}(\hat{\beta}_0, \hat{\boldsymbol{\beta}}) [\mathbf{1} \ \mathbf{X}] + \lambda_1 \boldsymbol{\Omega} + \lambda_2 \tilde{\boldsymbol{\Lambda}})^{-1} [\mathbf{1} \ \mathbf{X}]^\top [\mathbf{W}(\hat{\beta}_0, \hat{\boldsymbol{\beta}})]^{1/2} \mathbf{z}.$$

## 5.6 Standard errors

The heuristic of the previous section also turns out to be useful to obtain approximate standard errors for the estimated coefficients. For linear models, we have

$$(\hat{\beta}_0 \ \hat{\boldsymbol{\beta}}^\top)^\top \approx \underbrace{([\mathbf{1} \ \mathbf{X}]^\top [\mathbf{1} \ \mathbf{X}] + \lambda_1 \boldsymbol{\Omega} + \lambda_2 \tilde{\boldsymbol{\Lambda}})^{-1} [\mathbf{1} \ \mathbf{X}]^\top \mathbf{y}}_{\boldsymbol{\Gamma}},$$

concluding that one may use standard errors

$$\text{se}(\hat{\beta}_j) = \hat{\sigma} \sqrt{(\boldsymbol{\Gamma} [\mathbf{1} \ \mathbf{X}]^\top [\mathbf{1} \ \mathbf{X}] \boldsymbol{\Gamma})_{jj}}, \quad j = 0, \dots, p,$$

where  $\hat{\sigma}$  denotes an estimator of the standard deviation of the error terms. Likewise for generalized linear models, we have the approximation

$$(\hat{\beta}_0 \ \hat{\boldsymbol{\beta}}^\top)^\top \approx \underbrace{([\mathbf{1} \ \mathbf{X}]^\top \mathbf{W}(\hat{\beta}_0, \hat{\boldsymbol{\beta}}) [\mathbf{1} \ \mathbf{X}] + \lambda_1 \boldsymbol{\Omega} + \lambda_2 \tilde{\boldsymbol{\Lambda}})^{-1} [\mathbf{1} \ \mathbf{X}]^\top \mathbf{W}(\hat{\beta}_0, \hat{\boldsymbol{\beta}}) \mathbf{z}}_{\boldsymbol{\Upsilon}(\hat{\beta}_0, \hat{\boldsymbol{\beta}})},$$

and consequently

$$\text{se}(\hat{\beta}_j) = \hat{\phi} \sqrt{\{\boldsymbol{\Upsilon}(\hat{\beta}_0, \hat{\boldsymbol{\beta}}) [\mathbf{1} \ \mathbf{X}]^\top \mathbf{W}(\hat{\beta}_0, \hat{\boldsymbol{\beta}}) [\mathbf{1} \ \mathbf{X}] \boldsymbol{\Upsilon}(\hat{\beta}_0, \hat{\boldsymbol{\beta}})\}_{jj}}, \quad j = 0, \dots, p.$$

Note that this approach is not adequate for those  $\hat{\beta}_j$  equal to zero, since it always yields  $\widehat{\text{se}}(\hat{\beta}_j) = 0$ . An alternative is the bootstrap, whose validity for computing standard errors in lasso-type estimation procedures has, up to our knowledge, not yet been studied.

## 5.7 Determining the hyperparameters

We propose to resort to the standard technique of cross-validation, though it is computationally expensive. For each value on a sufficiently fine grid of  $(\lambda_1, \lambda_2)$ -values, we compute the cross-validated loss

$$\text{CV}(\lambda_1, \lambda_2) = \sum_{l=1}^k \sum_{i: (\mathbf{x}_i, y_i) \in S_l} L(y_i, f(\mathbf{x}_i; \hat{\beta}_0^{S_{-l}}, \hat{\beta}^{S_{-l}}))$$

by randomly dividing  $S$  into cross-validation folds  $S_1, \dots, S_k$  of roughly equal size, defining  $S_{-l} = S \setminus S_l$ ,  $l = 1, \dots, k$ , and denoting by  $\hat{\beta}_0^{S_{-l}}, \hat{\beta}^{S_{-l}}$  the structured elastic net estimates using the sample  $S_l$ . Alternatively, one may compute a model selection criterion such as the GCV,  $C_p$ , AIC on the basis of the formula of degrees of freedom given in Section 5.5.

## 6 Data Analysis

### 6.1 One-dimensional signal regression

In one-dimensional signal regression, as described, e.g., in [Frank and Friedman \(1993\)](#), one aims at the prediction of a response given a sampled signal  $\mathbf{x}^\top = (x(t))_{t=1}^T$ , where the indices  $t = 1, \dots, T$ , refer to different ordered sampling points. For a sample  $S = \{(\{x_1(t)\}_{t=1}^T, y_1), \dots, (\{x_n(t)\}_{t=1}^T, y_n)\}$  of pairs consisting of sampled signals and responses, we consider prediction models of the form

$$\hat{y}_i = \zeta \left( \hat{\beta}_0 + \sum_{t=1}^T x_i(t) \hat{\beta}(t) \right), \quad i = 1, \dots, n.$$

#### 6.1.1 Simulation study

Similarly to [Tutz and Gertheiss \(2009\)](#), we simulate signals  $x(t)$ ,  $t = 1, \dots, T$ ,  $T = 100$ , according to

$$\begin{aligned} \{x(t)\} &\sim \sum_{k=1}^5 b_k \sin(t\pi(5 - b_k)/50 - m_k) + \tau(t), \\ \{b_k\} &\sim U(0; 5), \quad \{m_k\} \sim U(0; 2\pi), \quad \{\tau(t)\} \sim N(0, 0.25), \end{aligned}$$

with  $U(a; b)$  denoting the uniform distribution on the interval  $(a; b)$ . For the coefficient function  $\beta^*(t)$ ,  $t = 1, \dots, T$ , we examine two cases. In the first case, referred to as 'bump setting', we use

$$\beta^*(t) = \begin{cases} -\{(30 - t)^2 + 100\} / 200 & t = 21, \dots, 39, \\ \{(70 - t)^2 - 100\} / 200 & t = 61, \dots, 80, \\ 0 & \text{otherwise.} \end{cases}$$

In the second case, referred to as 'block setting',

$$\beta^* = (\underbrace{0, \dots, 0}_{20 \text{ times}}, \underbrace{0.5, \dots, 0.5}_{10 \text{ times}}, \underbrace{1, \dots, 1}_{10 \text{ times}}, \underbrace{0.5, \dots, 0.5}_{10 \text{ times}}, \underbrace{0.25, \dots, 0.25}_{10 \text{ times}}, \underbrace{0, \dots, 0}_{40 \text{ times}})^\top.$$

The form of the signals and coefficient functions are displayed in [Figure 6](#).

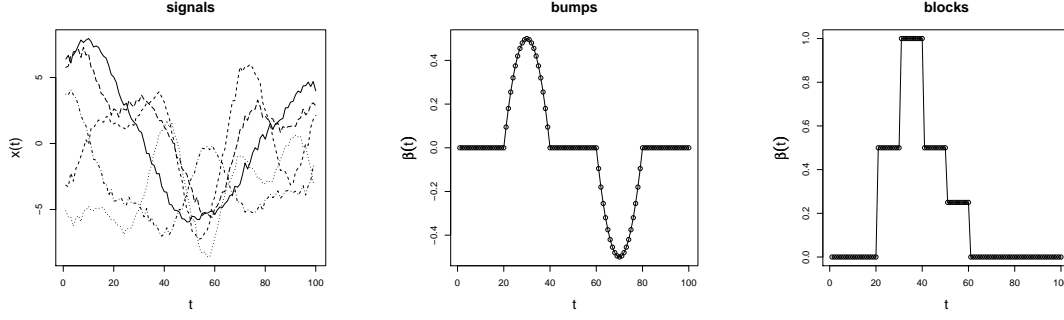


Figure 6: The setting of the simulation study. A collection of five signals (left panel), the coefficient functions for 'bump'- (middle panel) and 'block' setting (right panel), respectively.

For both settings, data are simulated according to

$$y = \sum_{t=1}^T x(t)\beta^*(t) + \epsilon, \quad \epsilon \sim N(0, 5),$$

For each out of 50 iterations, we simulate  $i = 1, \dots, 500$  i.i.d. realizations and divide them into three parts: a training set of size 200, a validation set of size 100 and a test set of size 200. Hyperparameters of the methods listed below are optimized by means of the validation set. As performance measures, we compute the absolute distance  $L^1(\hat{\beta}, \beta) = \|\hat{\beta} - \beta\|_1$  of true- and estimated coefficients and the mean squared prediction error on the test set. For methods with built-in feature selection, we additionally evaluate the goodness of selection in terms of sensitivity and specificity, defined by

$$\text{sensitivity}(\hat{\beta}, \beta^*) = \frac{|\{t : \hat{\beta}(t) \neq 0\} \cap \{t : \beta^*(t) \neq 0\}|}{|\{t : \beta^*(t) \neq 0\}|},$$

$$\text{specificity}(\hat{\beta}, \beta^*) = \frac{|\{t : \hat{\beta}(t) = 0\} \cap \{t : \beta^*(t) = 0\}|}{|\{t : \beta^*(t) = 0\}|}.$$

For each of the two setups, the simulation is repeated 50 times. The following methods are compared:

- ridge regression,
- generalized ridge regression with a first difference penalty,
- P-splines according to [Eilers and Marx \(1999\)](#),
- lasso,
- fused lasso,
- elastic net,
- structured elastic net with a first difference penalty,
- adaptive structured elastic net, where the weights  $\{\omega(t)\}$  are chosen according to the ridge estimator of the same iteration as  $\omega(t) = 1/|\hat{\beta}^{\text{ridge}}(t)|$ .

Performance measures are averaged over 50 iterations and displayed in [Table 1](#) (bump setting) and [Table 2](#) (block setting), respectively.

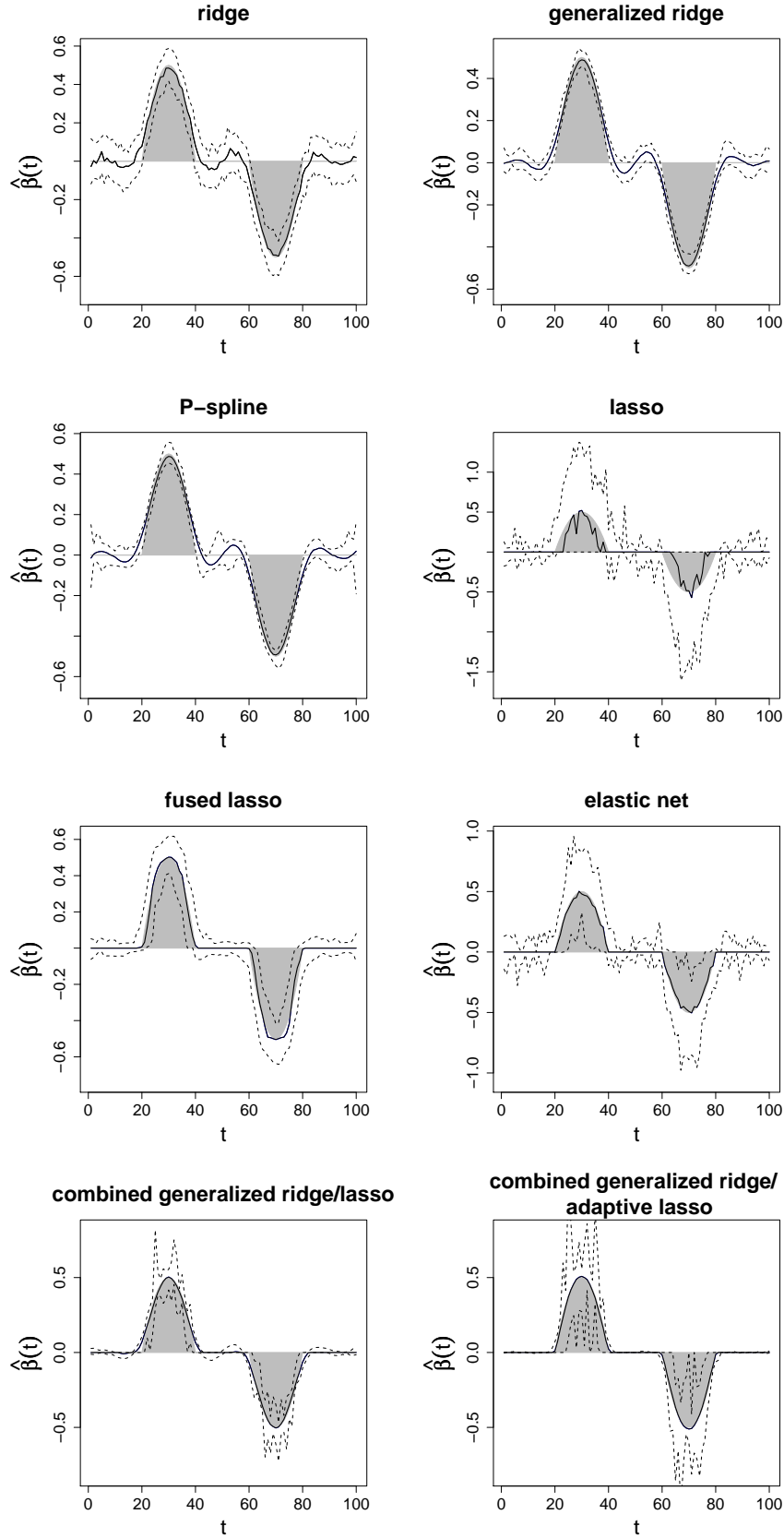


Figure 7: Estimated coefficient functions for the bump setting. The pointwise median curve over 50 iterations is represented by a solid line, pointwise 0.05- and 0.95-quantiles are drawn in dashed lines.

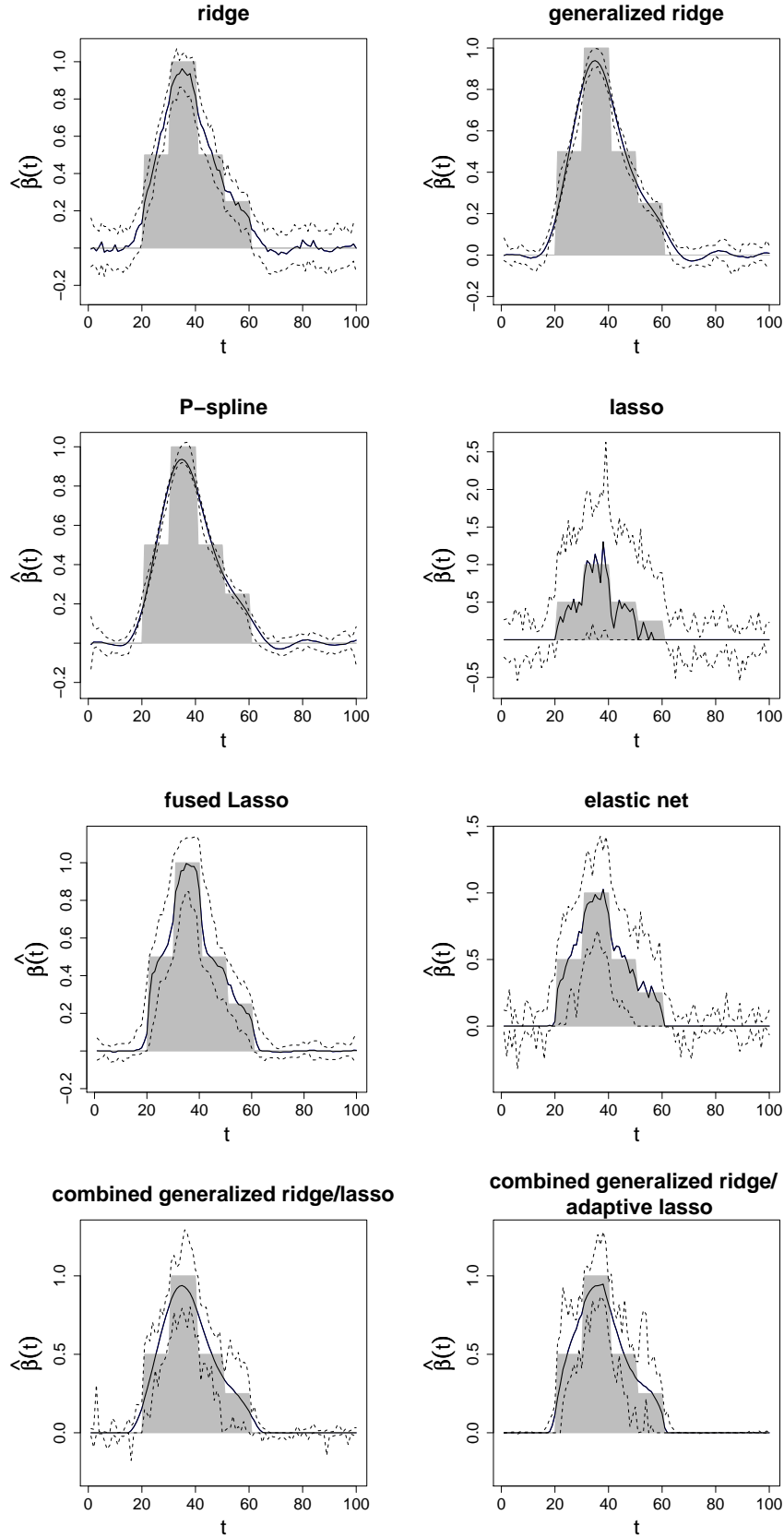


Figure 8: Estimated coefficient functions for the block setting. The pointwise median curve over 50 iterations is represented by a solid line, pointwise 0.05- and 0.95-quantiles are drawn in dashed lines.

method	$L^1(\widehat{\beta}, \beta^*)$	PE	sensitivity	specificity
ridge	0.249 ( $5.9 \times 10^{-4}$ )	5.35 (0.078)		
g.ridge	0.238 ( $9.9 \times 10^{-4}$ )	5.32 (0.076)		
P-spline	0.241 ( $16.0 \times 10^{-4}$ )	5.30 (0.077)		
lasso	0.271 ( $23.9 \times 10^{-4}$ )	5.72 (0.079)	0.62 ( $8.9 \times 10^{-3}$ )	0.65 (0.016)
fused lasso	0.235 ( $7.2 \times 10^{-4}$ )	5.30 (0.075)	0.96 ( $5.5 \times 10^{-3}$ )	0.51 (0.010)
enet	0.246 ( $29.9 \times 10^{-4}$ )	5.46 (0.081)	0.93 (0.013)	0.69 (0.032)
s.enet	<b>0.232</b> ( $7.6 \times 10^{-4}$ )	5.30 (0.078)	<b>0.98</b> ( $7.8 \times 10^{-3}$ )	0.59 (0.029)
ada.s.enet	<b>0.232</b> ( $15.0 \times 10^{-4}$ )	<b>5.25</b> (0.075)	0.91 ( $21.0 \times 10^{-3}$ )	<b>0.82</b> (0.020)

Table 1: Results for the bump setting, averaged over 50 simulations. We make use of the following abbreviations: 'PE' for 'mean squared prediction error', 'g. ridge' for 'generalized ridge', 'enet' for 'elastic net', 's.enet' for 'structured elastic net' and 'ada.s.enet' for 'adaptive structured elastic net'. Standard errors are given in parentheses. For each column, the best performance is emphasized in boldface.

For the bump setting, Figure 7 shows that at least the crude shape of the coefficient function is estimated by all compared methods in a satisfactory way, except for the lasso. Due to a favourable signal-to-noise ratio, even simplistic approaches such as ridge- or generalized ridge regression show competitive performance with respect to prediction of future observations. In pure numbers, the estimation of  $\beta^*$  is satisfactory as well. However, the lack of sparsity results into 'noise fitting' for those parts where  $\beta^*(t)$  is zero. For the two settings examined here, the P-spline approach does not improve over generalized ridge regression, because the two coefficient functions are not overly smooth. The elastic net considerably improves over the lasso and visually also over ridge regression, but it lacks smoothness. Its numerical inferiority to ridge regression results from double shrinkage as discussed in Subsection 3.5. The performance of the structured elastic net is not fully satisfactory. In particular, at the changepoints from zero- to nonzero parts, there is a tendency to widen unnecessarily the support of the nonzero sections. This shortcoming is removed by the adaptive structured elastic net, thereby confirming the theoretical result concerning selection consistency. This quality seems to be supported by the eminent performance with respect to sensitivity and specificity. The success of the adaptive strategy is also founded on the good performance of the ridge estimator providing the component-specific weights  $\omega(t)$ . The block setting is actually tailored to the fused lasso, whose output are piecewise constant coefficient functions. Nevertheless, it is not optimal, as the shrinkage of the  $\ell^1$ -penalty acts on all coefficients, including those different from zero. As a result, the fused lasso is outperformed by the adaptive structured elastic net with respect to prediction, though the structure part is seen to be not fully appropriate in the block setting. As opposed to the bump setting, fitting the block function seems to be much more difficult to accomplish in general.



method	$L^1(\hat{\beta}, \beta^*)$	PE	sensitivity	specificity
ridge	0.082 ( $3.4 \times 10^{-3}$ )	5.41 (0.080)		
g.ridge	0.064 ( $1.9 \times 10^{-3}$ )	5.35 (0.078)		
P-spline	0.065 ( $1.9 \times 10^{-3}$ )	5.34 (0.077)		
lasso	0.207 ( $3.6 \times 10^{-3}$ )	6.12 (0.089)	0.73 ( $7.5 \times 10^{-3}$ )	0.62 (0.014)
fused lasso	<b>0.058</b> ( $1.9 \times 10^{-3}$ )	5.34 (0.076)	<b>0.99</b> ( $7 \times 10^{-4}$ )	0.51 (0.009)
enet	0.094 ( $5.0 \times 10^{-3}$ )	5.47 (0.072)	0.95 ( $6.4 \times 10^{-3}$ )	0.73 (0.083)
s.enet	0.070 ( $5.0 \times 10^{-3}$ )	5.38 (0.080)	<b>0.99</b> ( $3.3 \times 10^{-3}$ )	0.60 (0.027)
ada.s.enet	0.061 ( $3.2 \times 10^{-3}$ )	<b>5.32</b> (0.69)	0.97 ( $8.0 \times 10^{-3}$ )	<b>0.83</b> (0.018)

Table 2: Results for the block setting, averaged over 50 simulations. For annotation, see Table 1.

### 6.1.2 Accelerometer data

The 'Sylvia Lawry Centre for Multiple Sclerosis Research e.V.', Munich, kindly provided us with two accelerometer records of two healthy female persons, aged between 20 and 30. They were equipped with a belt containing an accelerometer integrated into the belt buckle before walking several minutes on a flat surface at a moderate speed. The output are triaxial (vertical, horizontal, lateral) acceleration measurements at roughly 25,000 sampling points per person. Following [Daumer et al. \(2007\)](#) human gait, if defined as temporal evolution of three-dimensional accelerations of the center of mass of the body, is supposed to be a quasi-periodic process. Every period defines one gait cycle / double step, which starts with the heel strike and ends with the heel strike of the same foot. A single step ends with the heel strike of the other foot. Therefore, a double step can be seen as natural unit. As consequence, decomposition of the raw signal into pieces, each representing one double step, is an integral part of data preprocessing, not described in further detail here. Overall, we extract  $i = 1, \dots, n = 406$  double steps, 242 from person B ( $y = 0$ ) and 164 from person A ( $y = 1$ ), ending up with a sample  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ , where each  $\mathbf{x}_i = (x_i(t))$ ,  $t = 1, \dots, T = 102$  stores the observed vertical acceleration within double step  $i$ ,  $i = 1, \dots, n$ . For simplicity, we neglect the dependence of consecutive double steps within the same person and treat them as independent realizations. Horizontal- and lateral acceleration are not considered, since they do not carry information relevant to our prediction problem. We aim at the prediction of the person (A or B) given a double step pattern, and additionally at the detection of parts of the signal apt for discriminating between the two persons. We randomly divide the complete sample into a learning set of size 300 and a test set of size 106, and subsequently carry out logistic regression on the training set, using the structured elastic net with a squared first difference penalty. Hyperparameters are determined by ten-fold cross-validation, and the resulting logistic regression model is used to obtain predictions for the test set. The fused lasso with the hinge loss of support vector machines is used as competitor. A collection of results is assembled in Figures 9, 10 and Table 3, from which one concludes that classification is an easy task, since (nearly) perfect misclassification rates

on the test set are achieved. Concerning feature selection, the results of the structured elastic net are comparable to those of the fused lasso.

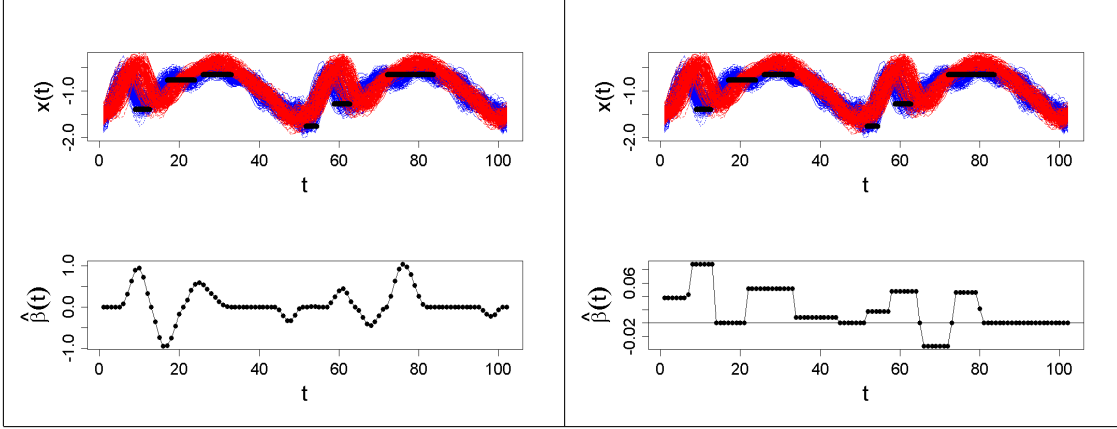


Figure 9: Coefficient functions for structured elastic net-regularized logistic regression (left panel) and the fused lasso support vector machine (right panel). Within each panel, the upper panel displays the overlaid double step patterns of the complete sample (406 double steps). Colours refer to persons, and visual differences between the two classes have been drawn manually as black bars.

fused lasso				
$t_1$	$t_2$	test error	# support vectors	# nonzero coefficients
2.5	0.5	0	124	46
structured elastic net				
$\lambda_1$	$\lambda_2$	test error	degrees of freedom	# nonzero coefficients
1.5	5	1	7.85	61

Table 3: Results of step classification for the fused lasso support vector machine and structured elastic net-regularized logistic regression. For the fused lasso,  $t_1$  denotes the bound imposed on the 1-norm of  $\beta$  corresponding to  $\lambda_1$ , and  $t_2$  denotes the bound on the absolute differences  $\sum_{t=2}^T |\beta(t) - \beta(t-1)|$ . For estimating the degrees of freedom, we make use of the heuristic suggested in Section 5.

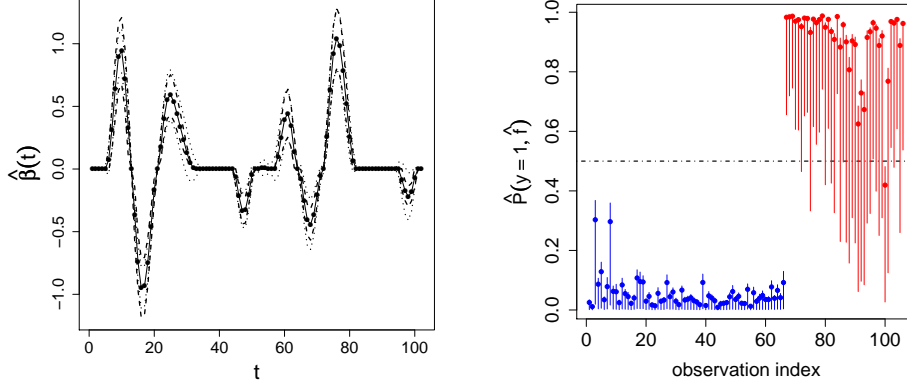


Figure 10: Left panel: the coefficient function of structured elastic net logistic regression (black points) and pointwise  $\pm 2$  standard error curves, once computed according to heuristic proposed in Section 5 (dashed lines) and once using 100 bootstrap iterations with fixed hyperparameters (dotted lines). Right panel: predicted probabilities for class 1, evaluated on the test set, indicated by filled points. The bars quantify uncertainty: lower- and upper end correspond to 0.1- and 0.9-quantile, respectively, computed from the bootstrap.

## 6.2 Surface fitting

Figure 11 depicts the surface to be fitted on a  $20 \times 20$  grid. The surface can be represented by a discrete function  $\beta^*(t, u)$ ,  $t, u = 1, \dots, 20$ . It consists of three non-overlapping truncated Gaussians of different shape and one plateau function. We have

$$\begin{aligned} \beta^*(t, u) &= B(t, u) + G_1(t, u) + G_2(t, u) + G_3(t, u), \\ B(t, u) &= \frac{1}{2} I(t \in \{10, 11, 12\}, u \in \{3, 4\}), \\ G_1(t, u) &= \max \left\{ 0, \exp \left( -(t-3 \ u-8) \begin{pmatrix} 3 & 0 \\ 0 & 0.25 \end{pmatrix} \begin{pmatrix} t-3 \\ u-8 \end{pmatrix} \right) - 0.2 \right\}, \\ G_2(t, u) &= \max \left\{ 0, \exp \left( -(t-7 \ u-17) \begin{pmatrix} 0.75 & 0 \\ 0 & 0.75 \end{pmatrix} \begin{pmatrix} t-7 \\ u-17 \end{pmatrix} \right) - 0.2 \right\}, \\ G_3(t, u) &= \max \left\{ 0, \exp \left( -(t-15 \ u-14) \begin{pmatrix} 0.5 & -0.25 \\ -0.25 & 0.5 \end{pmatrix} \begin{pmatrix} t-15 \\ u-14 \end{pmatrix} \right) - 0.2 \right\} \end{aligned} \quad (6.1)$$

Similarly to the simulation study in Subsection 6.1.1, we simulate a noisy version of the surface according to

$$y(t, u) = \beta^*(t, u) + \epsilon(t, u), \quad \{\epsilon(t, u)\} \stackrel{\text{i.i.d.}}{\sim} N(0, 0.25^2), \quad t, u = 1, \dots, 20.$$

For each of the 50 runs, we simulate two instances of  $y(t, u)$ . The first one is used for training and the second one for hyperparameter tuning. The mean squared error for estimating  $\beta^*$  is computed and averaged over 50 runs. Results are summarized in Figure 12 and Table 4. We compare ridge, generalized ridge with a difference penalty according to the grid structure, lasso, fused lasso with a total variation penalty along the grid, structured- and adaptive structured elastic net with the same difference penalty as for generalized ridge. The elastic net coincides - up to a constant scaling factor - with the lasso/soft thresholding in the orthogonal design case and is hence not considered.

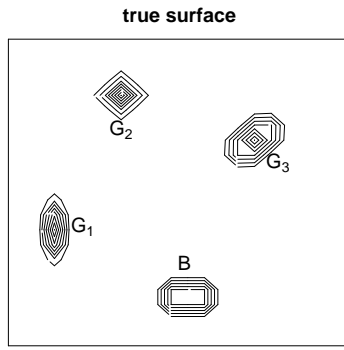


Figure 11: Contours of the surface according to Eq. (6.1).

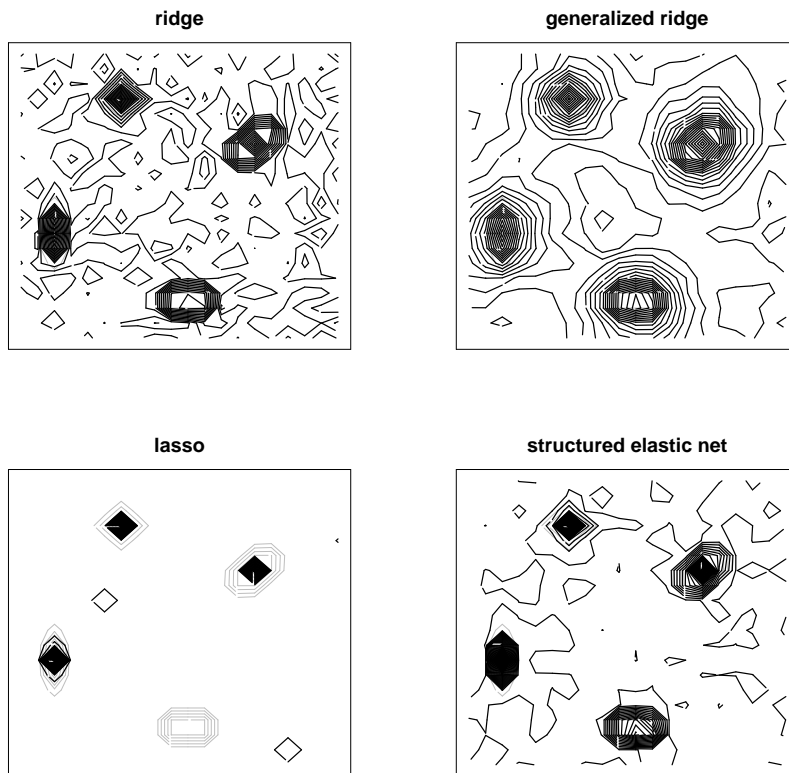


Figure 12: Contours of the estimated surfaces for four selected methods, averaged pointwise over 50 runs.

method	PE	$B$	$G_1$	$G_2$	$G_3$	zero
ridge	1.20 (0.01)	0.31	0.28	0.20	0.34	0.07
g.ridge	1.17 (0.04)	0.18	0.20	0.14	0.21	0.44
lasso	1.31 (0.01)	0.37	0.32	0.22	0.39	<b>0.01</b>
fused lasso	0.67 (0.02)	<b>0.14</b>	0.12	<b>0.08</b>	<b>0.15</b>	0.18
s.enet	0.88 (0.02)	0.22	0.16	0.12	0.23	0.18
ada.s.enet	<b>0.56</b> (0.02 )	0.15	<b>0.09</b>	<b>0.08</b>	0.18	0.06

Table 4: Results of the simulation, averaged over 50 iterations (standard errors in parentheses). The columns labeled  $B$ ,  $G_1$ ,  $G_2$ ,  $G_3$  and 'zero' contain the mean prediction error for the corresponding region of the surface. The abbreviations equal those in Table 1. The prediction error has been rescaled by 100.

### 6.3 Leukaemia cancer data

The dataset of [Golub et al. \(1999\)](#) constitutes one of the milestones in molecular classification of cancer. It consists of gene expression intensities for 7129 genes of 38 leukaemia patients, from which 27 are diagnosed acute lymphoblastic leukaemia (ALL), and the remaining patients are diagnosed acute myeloid leukaemia (AML). In addition, there is an independent test set of 34 samples. A major challenge in cancer research is the detection of pathways governing or influencing the rise and development of cancer. While this rationale is used to motivate the elastic net, it lacks the explicit integration of information on pathways, as available, e.g., in the KEGG database ([Kanehisa and Goto \(2000\)](#)). The latter represents metabolic pathways as graphs in which chemical reactions form the vertex set, and the edges are labeled by proteins/genes taking part in or catalyzing these reactions. [Li and Li \(2008\)](#) incorporate this information by employing what they term 'network-constrained regularization and variable selection', a special case of the structured elastic net, where the matrix  $\mathbf{A}$  is chosen as the *normalized* graph Laplacian ([Chung \(1997\)](#)), a rescaled version of the combinatorial Laplacian introduced in Section 2.2. Our approach for cancer class prediction making use of knowledge about KEGG pathways operates in a similar manner. We are able to match 2761 out of 7129 genes in the KEGG database, i.e. these 2761 genes occur in at least one KEGG pathway. Next, we construct an unweighted association graph by connecting two genes if and only if they share at least one common pathway. The resulting graph consists of eight connected components, from which the largest has a vertex set of size 2688. Following [Li and Li \(2008\)](#), we compute the normalized graph Laplacian, whose rank equals the number of vertices minus the number of connected components. The normalized graph Laplacian constitutes the matrix  $\mathbf{A}$  of structured net-regularized logistic regression model involving all 7129 genes. Since we are given only 38 observations, it follows from the augmented data representation that the structured elastic net degenerates in our situation. Concerning the aim of our analysis, we are primarily interested in the amount the structured elastic net differs from the elastic net in terms of variable selection. We do not observe a notable difference concerning prediction error on the test set, which is not surprising for three reasons. Firstly, prediction is not a hard task for the dataset at hand, because even simple approaches such as the nearest

enet		$\lambda_2 \xrightarrow{+}$		
# selected	12 (26)	19 (41)	27 (56)	64 (156)
$\ \hat{\beta}\ _1$	4.49 (9.27)	4.64 (9.29)	4.65 (9.34)	4.66 (9.90)
s.enet				
# selected	19 (44)	23 (73)	27 (91)	36 (135)
$\ \hat{\beta}\ _1$	0.016 (9.18)	0.004 (9.39)	0.003 (9.60)	0.003 (10.02)

Table 5: The table displays the number of selected pathway genes for the elastic net and the structured elastic net and the 1-norm of the corresponding coefficient subvector for selected values of  $\lambda_2$ , increasing from left to right and  $\lambda_1$  kept fixed. The numbers in parentheses refer to all genes, including the non-pathway genes.

shrunk centroid classifier (Tibshirani et al. (2003)) yield good results. Secondly, as it becomes clear from Table 5, the genes *not* mapped to any pathway are the relevant ones. Finally, we are sufficiently self-critical to admit that the 38 bits of information contained in the response are probably not enough to judge accurately the influence of pathways. Concerning selection of variables occurring in pathways, we observe striking differences when comparing the elastic net and its structured counterpart. Table 5 shows that for large  $\lambda_2$ , the importance of pathway-based predictors tends to zero for the structured elastic net, while it remains roughly constant for the elastic net. As shown in Figures 13 and 14, the number of pathway genes increases for both methods when increasing  $\lambda_2$ , and the increase is stronger for the elastic net. To judge the importance of pathways, we additionally apply subset selection based on gene set enrichment analysis (GSA, Efron and Tibshirani (2007)), in which the members of the same pathways are treated as one group and the non-pathway genes are treated as groups consisting of one member only. We include a group as a whole in a ridge-regularized logistic regression model if it is given a high score in GSA. The latter is applied repeatedly in ten-fold cross-validation on the training set to select the optimal number of groups. The results of this procedure confirm the observations made for the structured elastic net, because we end up with a competitive prediction model that only includes non-pathway genes.



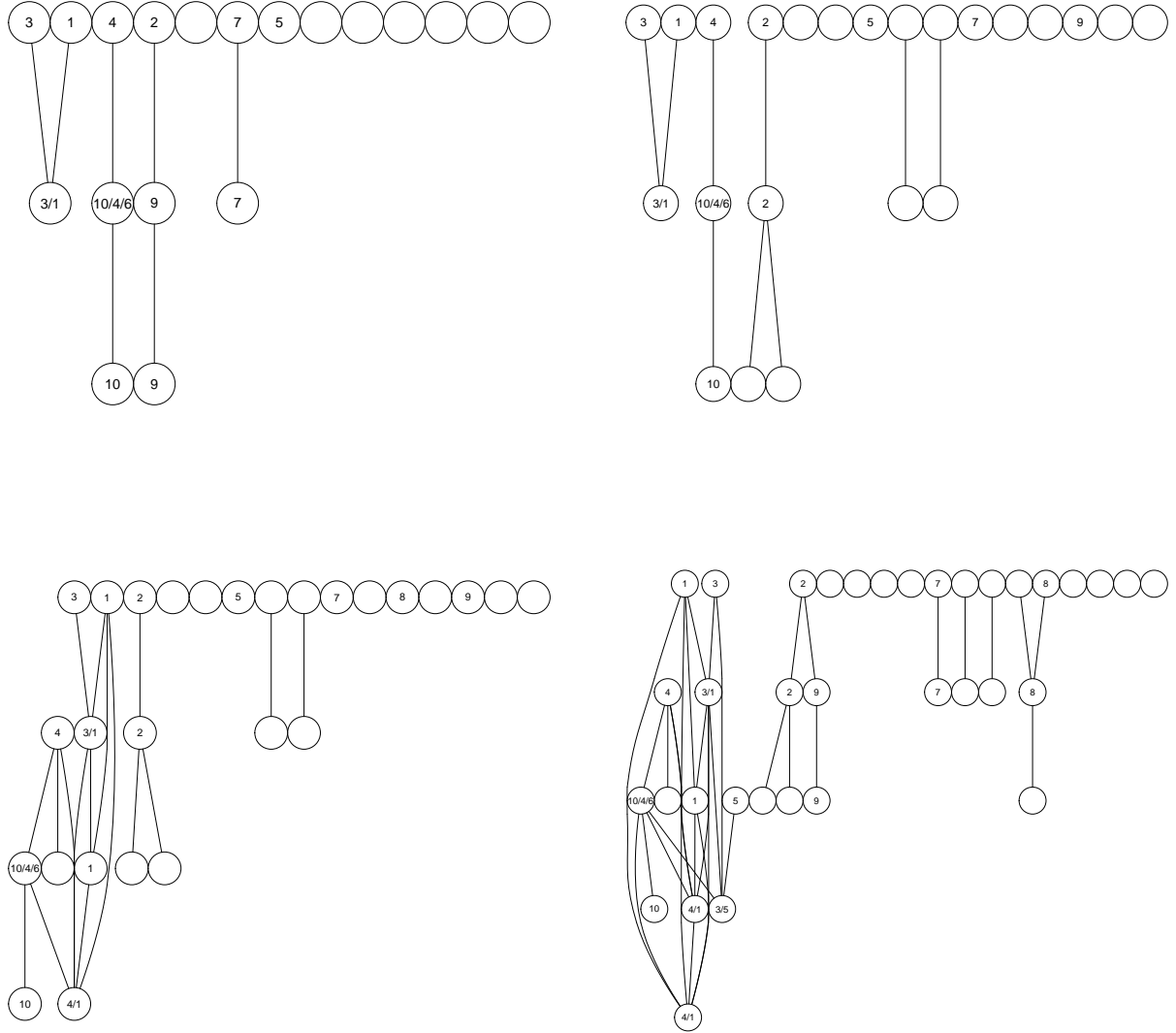


Figure 14: Selected pathway genes of the structured elastic net in the situation of Table 5, with  $\lambda_2$  increasing from top to bottom and left to right. The numerals refer to frequently occurring pathways: 1 - 'hematopoietic cell lineage', 2 - 'epithelial cell signaling', 3 - 'cytokine-cytokine receptor interaction', 4 - 'ECM receptor interaction', 5 - 'cell cycle', 6 - 'bladder cancer', 7 - 'arachidonic acid metabolism', 8 - 'Tryptophan metabolism', 9 - 'calcium signaling pathway', 10 - 'focal adhesion'.



## 7 Discussion

The structured elastic net is proposed as a procedure for coefficient selection and -smoothing. We have established a general notion of structured features, for which the structured elastic net is able to take advantage of prior knowledge as opposed to the lasso and the elastic net, which are both purely data-driven. The structured elastic net may also be regarded as a computationally more convenient alternative to the fused lasso. Conceptually, generalizing the fused lasso by computing the total variation of the coefficients along a graph is straightforward. However, due to the non-differentiability of the structure part of the fused lasso, computation may be intractable even for moderately sized graphs.

Turning to the drawbacks of the structured elastic net, we have outlined in Section 5 that model selection and computation of standard errors and in turn the quantification of uncertainty is notoriously difficult. A Bayesian approach promises to be superior in this regard. The lasso can be treated within a Bayesian inference framework (Park and Casella (2008)), while the quadratic part of the structured elastic net regularizer is already motivated from a Bayesian perspective in this paper.

With regard to possible directions of future research, we consider to study the structured elastic net in combination with other loss functions, e.g. the hinge loss of support vector machines or the check loss for quantile regression. Concerning the application to genomic data, we work on an extension to survival data, in particular to the Cox proportional hazards model. The asymptotic analysis in this paper is basic in the sense that it is bound to strong assumptions, and the role of the structure part of the regularizer and its interplay with the true coefficient vector is not well understood yet, leaving some room for more profound investigations.

## Acknowledgments

Martin Slawski was partly supported by the Porticus Foundation for Clinical Medicine and Bioinformatics. We thank the Sylvia Lawry Centre for its support with the accelerometer data example, in particular Martin Daumer for numerous discussions, Christine Gerges and Kathrin Thaler for producing the data. We thank Jelle Goeman for one helpful discussion about his algorithm and making his code publicly available as R package. We are grateful to Angelika van der Linde and Daniel Sabanés-Bové for pointing us to several errors and typos in earlier drafts.

## A Proofs

### A.1 Proof of Proposition 1

Using the well-known expression for the gradient  $\frac{\partial L}{\partial \beta}$  in generalized linear models, the Karush-Kuhn-Tucker (KKT) conditions imply that

$$\begin{aligned} -\mathbf{X}_1^\top (\mathbf{y} - \hat{\boldsymbol{\mu}}) + \lambda_2(\hat{\beta}_1 + s\hat{\beta}_2) + \lambda_1 \text{sign}(\hat{\beta}_1) &= 0, \\ -\mathbf{X}_2^\top (\mathbf{y} - \hat{\boldsymbol{\mu}}) + \lambda_2(\hat{\beta}_2 + s\hat{\beta}_1) + \lambda_1 \text{sign}(\hat{\beta}_2) &= 0, \end{aligned}$$

where  $\hat{\boldsymbol{\mu}} = \hat{\mathbb{E}}[\mathbf{y} | \mathbf{X}_1, \mathbf{X}_2]$ . Adding the second equation to the first equation multiplied by  $s$  yields

$$-(s\mathbf{X}_1 + \mathbf{X}_2)^\top (\mathbf{y} - \hat{\boldsymbol{\mu}}) + 2\lambda_2(\hat{\beta}_2 + s\hat{\beta}_1) = 0,$$

implying

$$\begin{aligned}
|\hat{\beta}_1 + s\hat{\beta}_2| &= \frac{1}{2\lambda_2} |(s\mathbf{X}_1 + \mathbf{X}_2)^\top (\mathbf{y} - \hat{\boldsymbol{\mu}})| \\
&\leq \frac{1}{2\lambda_2} \|s\mathbf{X}_1 + \mathbf{X}_2\| \|\mathbf{y} - \hat{\boldsymbol{\mu}}\| \\
&\leq \frac{1}{2\lambda_2} \sqrt{2(1+s\rho)} \|\mathbf{y}\|,
\end{aligned}$$

noting that  $\|s\mathbf{X}_1 + \mathbf{X}_2\| = (\|\mathbf{X}_1\|^2 + \|\mathbf{X}_2\|^2 + 2s\langle \mathbf{X}_1, \mathbf{X}_2 \rangle)^{1/2} = (2(1+s\rho))^{1/2}$ , since  $\mathbf{X}_1$  and  $\mathbf{X}_2$  are standardized.

## A.2 Proof of Proposition 2

Consider the case  $s = -1$  and define another set of coefficients by

$$\tilde{\beta}_0 = \hat{\beta}_0, \quad \tilde{\beta}_1 = \tilde{\beta}_2 = \frac{1}{2}(\hat{\beta}_1 + \hat{\beta}_2).$$

First observe that

$$\begin{aligned}
&\sum_{i=1}^n L(y_i, f(\mathbf{x}_i; \tilde{\beta}_0, \tilde{\boldsymbol{\beta}})) + \frac{\lambda_2}{2} (\tilde{\beta}_1 - \tilde{\beta}_2)^2 + \lambda_1 \|\tilde{\boldsymbol{\beta}}\|_1 \\
&- \sum_{i=1}^n L(y_i, f(\mathbf{x}_i; \hat{\beta}_0, \hat{\boldsymbol{\beta}})) - \frac{\lambda_2}{2} (\hat{\beta}_1 - \hat{\beta}_2)^2 - \lambda_1 \|\hat{\boldsymbol{\beta}}\|_1 \geq 0.
\end{aligned} \tag{A.1}$$

For the part involving the loss function, we obtain

$$\begin{aligned}
&\sum_{i=1}^n \left\{ L(y_i, f(\mathbf{x}_i; \tilde{\beta}_0, \tilde{\boldsymbol{\beta}})) - L(y_i, f(\mathbf{x}_i; \hat{\beta}_0, \hat{\boldsymbol{\beta}})) \right\} \\
&\leq \sum_{i=1}^n \left| L(y_i, f(\mathbf{x}_i; \tilde{\beta}_0, \tilde{\boldsymbol{\beta}})) - L(y_i, f(\mathbf{x}_i; \hat{\beta}_0, \hat{\boldsymbol{\beta}})) \right| \\
&\leq C \sum_{i=1}^n |\mathbf{x}_i^\top (\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}})| \\
&= C \sum_{i=1}^n \left| \frac{1}{2} (x_{i1} - x_{i2}) (\tilde{\beta}_2 - \hat{\beta}_1) \right| \\
&= \frac{C}{2} |\tilde{\beta}_2 - \hat{\beta}_1| \|\mathbf{X}_1 - \mathbf{X}_2\|_1.
\end{aligned} \tag{A.2}$$

For the  $\ell^1$ -regularizer, we have

$$\|\tilde{\boldsymbol{\beta}}\|_1 - \|\hat{\boldsymbol{\beta}}\|_1 = |\tilde{\beta}_1 + \tilde{\beta}_2| - |\hat{\beta}_1| + |\hat{\beta}_2| \leq 0. \tag{A.3}$$

For the quadratic regularizer one obtains the difference  $\frac{1}{2}(\hat{\beta}_1 - \hat{\beta}_2)^2$ . Combining this with (A.1) - (A.3), we obtain

$$\frac{C}{2} |\hat{\beta}_1 - \hat{\beta}_2| \|\mathbf{X}_1 - \mathbf{X}_2\|_1 - \frac{\lambda_2}{2} (\hat{\beta}_1 - \hat{\beta}_2)^2 \geq 0,$$

from which it follows that

$$\begin{aligned} |\hat{\beta}_1 - \hat{\beta}_2| &\leq \frac{C}{\lambda_2} \|\mathbf{X}_1 - \mathbf{X}_2\|_1 \\ &\leq \frac{C}{\lambda_2} \sqrt{n} \|\mathbf{X}_1 - \mathbf{X}_2\| \\ &\leq \frac{C}{\lambda_2} \sqrt{n2(1-\rho)}, \end{aligned}$$

as shown in Proposition 1. The case  $s = 1$  is obtained analogously by setting  $\tilde{\beta}_0 = \hat{\beta}_0$ ,  $\tilde{\beta}_1 = \frac{1}{2}(\hat{\beta}_1 - \hat{\beta}_2)$ ,  $\tilde{\beta}_2 = \frac{1}{2}(\hat{\beta}_2 - \hat{\beta}_1)$ .

### A.3 Auxiliary results

The asymptotic analysis relies on a more general theory of constrained M-estimation involving the notion of pointwise convergence in distribution of convex functions.

**Definition A. 1.** Let  $G_n$  be a sequence of lower semicontinuous convex random functions from  $\mathbb{R}^p$  to  $\mathbb{R} \cup \{\infty\}$ , let  $G$  be another such function and let  $\mathcal{D}$  be a countable dense set in  $\mathbb{R}^p$ . Then  $G_n$  converges pointwise in distribution to  $G$ , in signs  $G_n \xrightarrow{D} G$ , if for each finite subset  $\{\mathbf{u}_1, \dots, \mathbf{u}_k\} \subset \mathcal{D}$ ,  $(G_n(\mathbf{u}_1), \dots, G_n(\mathbf{u}_k))^\top \xrightarrow{D} (G(\mathbf{u}_1), \dots, G(\mathbf{u}_k))^\top$ .

**Theorem A. 1.** ([Geyer \(1996\)](#))

Let  $G_n, G$  be as in Definition A. 1 and let  $G_n \xrightarrow{D} G$ . Define  $\hat{\mathbf{u}}_n = \operatorname{argmin} G_n$  and  $\hat{\mathbf{u}} = \operatorname{argmin} G$ . If  $G$  has a unique minimizer, then  $\hat{\mathbf{u}}_n \xrightarrow{D} \hat{\mathbf{u}}$ .

By means of this preparation, we can prove Theorems 1-4.

### A.4 Proof of Theorem 1

Define the random function  $V_n(\mathbf{u})$  by

$$\begin{aligned} V_n(\mathbf{u}) &= \sum_{i=1}^n (\epsilon_i - \mathbf{u}^\top \mathbf{x}_i / \sqrt{n})^2 - \epsilon_i^2 \\ &\quad + \lambda_1^n \sum_{j=1}^p |\beta_j^* + u_j / \sqrt{n}| - |\beta_j^*| \\ &\quad + \lambda_2^n [(\boldsymbol{\beta}^* + \mathbf{u} / \sqrt{n})^\top \boldsymbol{\Lambda} (\boldsymbol{\beta}^* + \mathbf{u} / \sqrt{n}) - \boldsymbol{\beta}^{*\top} \boldsymbol{\Lambda} \boldsymbol{\beta}^*]. \end{aligned}$$

Observe that  $V_n(\mathbf{u})$  is minimized at  $\hat{\mathbf{u}}_n = \sqrt{n}(\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}^*)$ , because with  $\mathbf{u} = \sqrt{n}(\boldsymbol{\beta} - \boldsymbol{\beta}^*)$ ,

$$\hat{\mathbf{u}}_n = \operatorname{argmin}_{\mathbf{u}} V_n(\mathbf{u}) \Leftrightarrow \hat{\boldsymbol{\beta}}_n = \operatorname{argmin}_{\boldsymbol{\beta}} (Q_n(\boldsymbol{\beta}) - Q_n(\boldsymbol{\beta}^*)),$$

with

$$Q_n(\boldsymbol{\beta}) = \|\mathbf{y}_n - \mathbf{X}_n \boldsymbol{\beta}\|^2 + \lambda_1^n \|\boldsymbol{\beta}\|_1 + \lambda_2^n \boldsymbol{\beta}^\top \boldsymbol{\Lambda} \boldsymbol{\beta}. \quad (\text{A.4})$$

Considering the first term of  $V_n(\mathbf{u})$ , we have

$$\begin{aligned} \sum_{i=1}^n (\epsilon_i - \mathbf{u}^\top \mathbf{x}_i / \sqrt{n})^2 - \epsilon_i^2 &= \frac{\mathbf{u}^\top \mathbf{X}_n^\top \mathbf{X}_n \mathbf{u}}{n} - 2\mathbf{u}^\top \sqrt{n} \frac{\mathbf{X}_n^\top (\mathbf{y}_n - \mathbf{X}_n \boldsymbol{\beta}^*)}{n} \\ &= \frac{\mathbf{u}^\top \mathbf{X}_n^\top \mathbf{X}_n \mathbf{u}}{n} - 2\mathbf{u}^\top \frac{\mathbf{X}_n^\top \mathbf{X}_n}{n} \sqrt{n} \left( \left( \frac{\mathbf{X}_n^\top \mathbf{X}_n}{n} \right)^{-1} \frac{\mathbf{X}_n^\top \mathbf{y}_n}{n} - \boldsymbol{\beta}^* \right). \end{aligned}$$

The first term on the r.h.s. converges to  $\mathbf{u}^\top \mathbf{C} \mathbf{u}$  by condition (C.2). The asymptotic result for the ols estimator is

$$\sqrt{n} \left( \left( \frac{\mathbf{X}_n^\top \mathbf{X}_n}{n} \right)^{-1} \frac{\mathbf{X}_n^\top \mathbf{y}_n}{n} - \boldsymbol{\beta}^* \right) \xrightarrow{D} N(\mathbf{0}, \sigma^2 \mathbf{C}^{-1}),$$

hence the second term of the previous expression converges to  $\mathbf{w}$  in distribution. Invoking Slutsky's theorem, the first term of  $V_n$  converges to  $-2\mathbf{u}^\top \mathbf{w} + \mathbf{u}^\top \mathbf{C} \mathbf{u}$ , again in distribution. For the first penalty term in  $V_n$ , one has that

$$\begin{aligned} \lambda_1^n \sum_{j=1}^p |\beta_j^* + u_j/\sqrt{n}| - |\beta_j^*| &= \frac{\lambda_1^n}{\sqrt{n}} \sum_{j=1}^p \{ \text{sign}(\beta_j^* + u_j/\sqrt{n}) (\sqrt{n}\beta_j^* + u_j) - \text{sign}(\beta_j^*) \sqrt{n}\beta_j^* \} I(\beta_j^* \neq 0) \\ &+ \frac{\lambda_1^n}{\sqrt{n}} \sum_{j=1}^p |u_j| I(\beta_j^* = 0). \end{aligned}$$

Since  $\text{sign}(\beta_j^* + u/\sqrt{n}) \rightarrow \text{sign}(\beta_j^*)$  and  $\lambda_1^n/\sqrt{n} \rightarrow \lambda_1^0$  as  $n \rightarrow \infty$ , the expression converges to the second line in the definition of  $V$ . Finally,

$$\lim_{n \rightarrow \infty} \lambda_2^n \{ (\boldsymbol{\beta}^* + \mathbf{u}/\sqrt{n})^\top \boldsymbol{\Lambda} (\boldsymbol{\beta}^* + \mathbf{u}/\sqrt{n}) - \boldsymbol{\beta}^{*\top} \boldsymbol{\Lambda} \boldsymbol{\beta}^* \} = 2\lambda_2^0 \mathbf{u}^\top \boldsymbol{\Lambda} \boldsymbol{\beta}^*.$$

Applying Slutsky's theorem once again, it holds that  $V_n \xrightarrow{D} V$  in the sense of Definition A. 1. Since  $V_n$  is convex and  $V$  has a unique minimizer, we conclude from Theorem A. 1 that

$$\underset{\mathbf{u}}{\text{argmin}} V_n(\mathbf{u}) = \hat{\mathbf{u}}_n \xrightarrow{D} \underset{\mathbf{u}}{\text{argmin}} V(\mathbf{u}) = \sqrt{n}(\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}^*).$$

## A.5 Proof of Theorem 2

Define  $W_n(\mathbf{u})$  by

$$\begin{aligned} W_n(\mathbf{u}) &= L_n(\mathbf{u}) + \lambda_1^n \sum_{j=1}^p |\beta_j^* + u_j/\sqrt{n}| - |\beta_j^*| + \lambda_2^n [(\boldsymbol{\beta}^* + \mathbf{u}/\sqrt{n})^\top \boldsymbol{\Lambda} (\boldsymbol{\beta}^* + \mathbf{u}/\sqrt{n}) - \boldsymbol{\beta}^{*\top} \boldsymbol{\Lambda} \boldsymbol{\beta}^*], \\ L_n(\mathbf{u}) &= 2\phi^{-1} \sum_{i=1}^n b(\mathbf{x}_i^\top (\mathbf{u}/\sqrt{n} + \boldsymbol{\beta}^*)) - b(\mathbf{x}_i^\top \boldsymbol{\beta}^*) - y_i \mathbf{x}_i^\top \mathbf{u}/\sqrt{n}. \end{aligned}$$

Observe that  $\underset{\mathbf{u}}{\text{argmin}} W(\mathbf{u}) = \hat{\mathbf{u}}_n = \sqrt{n}(\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}^*)$  with  $\hat{\boldsymbol{\beta}}_n$  as in Theorem 2. We have the following second-order Taylor expansion of  $L_n(\mathbf{u})$  around  $\mathbf{u} = \mathbf{0}$ :

$$\begin{aligned} L_n(\mathbf{u}) &= 2\phi^{-1} \sum_{i=1}^n (y_i - b'(\mathbf{x}_i^\top \boldsymbol{\beta}^*)) \frac{\mathbf{u}^\top \mathbf{x}_i}{\sqrt{n}} \\ &+ \phi^{-1} \sum_{i=1}^n b''(\mathbf{x}_i^\top \boldsymbol{\beta}^*) \frac{\mathbf{u}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{u}}{n} \\ &+ R_n(\mathbf{u}), \end{aligned}$$

where the remainder is given by

$$R_n(\mathbf{u}) = \frac{1}{3n^{3/2}} \phi^{-1} \sum_{i=1}^n b'''(\mathbf{x}_i^\top \boldsymbol{\xi}) (\mathbf{x}_i^\top \mathbf{u})^3,$$

where  $\boldsymbol{\xi}$  is contained in the segment from  $\boldsymbol{\beta}^*$  to  $\boldsymbol{\beta}^* + \mathbf{u}/\sqrt{n}$ . Considering the first term of  $L_n(\mathbf{u})$ , standard properties of generalized linear models can be applied (McCullagh and Nelder (1989)):

$$\begin{aligned}\phi^{-1} \mathbf{u}^\top \mathbb{E}[\mathbf{x}_i(y_i - b'(\mathbf{x}_i^\top \boldsymbol{\beta}^*))] &= 0, \\ \text{var}[\phi^{-1}(y_i - b'(\mathbf{x}_i^\top \boldsymbol{\beta}^*))] &= \mathbf{u}^\top \frac{\mathbb{E}[\phi^{-1} b''(\mathbf{x}_i^\top \boldsymbol{\beta}^*) \mathbf{x}_i \mathbf{x}_i^\top]}{n} \mathbf{u} = \frac{\mathbf{u}^\top \mathcal{I} \mathbf{u}}{n}.\end{aligned}$$

Application of the central limit theorem yields that

$$\phi^{-1} \frac{\mathbf{u}^\top}{\sqrt{n}} \sum_{i=1}^n \mathbf{x}_i (b'(\mathbf{x}_i^\top \boldsymbol{\beta}^*) - y_i) \xrightarrow{D} \mathbf{u}^\top \mathbf{w}, \quad \mathbf{w} \sim N(\mathbf{0}, \mathcal{I}).$$

For the second term of  $L_n(\mathbf{u})$ , note that

$$\phi^{-1} \sum_{i=1}^n b''(\mathbf{x}_i^\top \boldsymbol{\beta}^*) \frac{\mathbf{u}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{u}}{n} \rightarrow \mathcal{I}.$$

Turning to the remainder,

$$3n^{1/2} \phi R_n(\mathbf{u}) = \frac{1}{n} \sum_{i=1}^n b'''(\mathbf{x}_i^\top \boldsymbol{\xi}) \leq \frac{1}{n} \sum_{i=1}^n M(\mathbf{x}_i) (\mathbf{x}_i^\top \mathbf{u})^3 \rightarrow \mathbb{E}[M(\mathbb{X}) |\mathbf{u}^\top \mathbb{X}|^3] < \infty$$

by condition (G.2), concluding that  $R_n(\mathbf{u}) = O_P(n^{-1/2})$ . The limiting behaviour of the regularizer in  $W_n(\mathbf{u})$  has already been studied in the proof of Theorem 1. As for the latter, Slutsky's theorem and Theorem A. 1 imply that  $\text{argmin } W_n(\mathbf{u}) \xrightarrow{D} \text{argmin } W(\mathbf{u})$ .

## A.6 Proof of Theorem 3

Define

$$\begin{aligned}\Psi_n(\mathbf{u}) &= \sum_{i=1}^n \left( \epsilon_i - \frac{\lambda_1^n}{n} \mathbf{x}_i^\top \mathbf{u} \right)^2 - \epsilon_i^2 \\ &\quad + \lambda_1^n \sum_{j=1}^p \left| \beta_j^* + \frac{\lambda_1^n}{n} \right| - |\beta_j^*| \\ &\quad + \lambda_2^n \left\{ \left( \boldsymbol{\beta}^* + \frac{\lambda_1^n}{n} \mathbf{u} \right)^\top \boldsymbol{\Lambda} \left( \boldsymbol{\beta}^* + \frac{\lambda_1^n}{n} \mathbf{u} \right) - \boldsymbol{\beta}^{*\top} \boldsymbol{\Lambda} \boldsymbol{\beta}^* \right\},\end{aligned}$$

and  $\Xi_n(\mathbf{u}) = \Psi_n(\mathbf{u}) \cdot (\lambda_1^n)^2/n$ . If  $\hat{\boldsymbol{\beta}}_n$  denotes the minimizer of  $Q_n(\boldsymbol{\beta})$  in Eq. (A.4), then  $\hat{\mathbf{u}}_n = \text{argmin } \Psi_n(\mathbf{u}) = \text{argmin } \Xi_n(\mathbf{u}) = n/\lambda_1^n (\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}^*)$ . Next, we consider the termwise limits within  $\Xi_n(\mathbf{u})$ . We have

$$\begin{aligned}\Xi_n(\mathbf{u}) &= \frac{1}{n} \mathbf{u}^\top \mathbf{X}_n^\top \mathbf{X}_n \mathbf{u} - 2 \frac{\boldsymbol{\epsilon}_n^\top \mathbf{X}_n}{\lambda_1^n} \\ &\quad + \sum_{j=1}^p \frac{n}{\lambda_1^n} \left( \left| \beta_j^* + \frac{\lambda_1^n}{n} \right| - |\beta_j^*| \right) \\ &\quad + \frac{n \lambda_2^n}{(\lambda_1^n)^2} \left\{ \left( \boldsymbol{\beta}^* + \frac{\lambda_1^n}{n} \mathbf{u} \right)^\top \boldsymbol{\Lambda} \left( \boldsymbol{\beta}^* + \frac{\lambda_1^n}{n} \mathbf{u} \right) - \boldsymbol{\beta}^{*\top} \boldsymbol{\Lambda} \boldsymbol{\beta}^* \right\}.\end{aligned}$$

The first term converges to  $\mathbf{u}^\top \mathbf{C} \mathbf{u}$ . For the second term, we have

$$\frac{\boldsymbol{\epsilon}_n^\top \mathbf{X}_n}{\lambda_1^n} = \underbrace{\sqrt{n} \left( \frac{\boldsymbol{\epsilon}_n^\top \mathbf{X}_n}{n} \right)}_{=O_P(1)} \frac{\sqrt{n}}{\lambda_1^n} = o_P(1).$$

The third terms converges to

$$P(\mathbf{u}) = \sum_{j=1}^p \text{sign}(\beta_j^*) I(\beta_j^* \neq 0) + |u_j| I(\beta_j^* = 0).$$

For the last term, we have the limit  $2R\mathbf{u}^\top \boldsymbol{\Lambda} \boldsymbol{\beta}^*$ . Defining  $\Xi(\mathbf{u}) = \mathbf{u}^\top \mathbf{C} \mathbf{u} + P(\mathbf{u}) + 2R\mathbf{u}^\top \boldsymbol{\Lambda} \boldsymbol{\beta}^*$ ,

$\Xi_n \xrightarrow{D} \Xi$  and by Theorem A. 1,  $\hat{\mathbf{u}}_n \xrightarrow{P} \argmin \Xi = \hat{\mathbf{u}}$ .

Since we have claimed selection consistency,

$$P(\hat{\beta}_{j,n} = 0) \rightarrow 1 \text{ for all } j \in A^c,$$

from which it follows that  $\hat{u}_j = 0 \forall j \in A^c$ . On the other hand, using the partitioning scheme of Eq. (4.3),  $\hat{\mathbf{u}}_A$  satisfies the equation

$$2\mathbf{C}_A \hat{\mathbf{u}}_A + 2\mathbf{C}_{AA^c} \mathbf{u}_{A^c} + \mathbf{s}_A + 2R\boldsymbol{\Lambda}_A \boldsymbol{\beta}_A^* = \mathbf{0}, \quad \mathbf{s}_A = (\text{sign}(\beta_j^*), j \in A)^\top.$$

$$\Rightarrow \hat{\mathbf{u}}_A = -\mathbf{C}_A^{-1} \left( \mathbf{C}_{AA^c} \mathbf{u}_{A^c} + \frac{\mathbf{s}_A}{2} + R\boldsymbol{\Lambda}_A \boldsymbol{\beta}_A^* \right).$$

Partial optimization of  $\Xi$  w.r.t. to  $\mathbf{u}_{A^c}$  amounts to the minimization of the following expression:

$$\mathbf{u}_{A^c}^\top \mathbf{C}_{A^c} \mathbf{u}_{A^c} + 2\mathbf{u}_{A^c}^\top \mathbf{C}_{A^c A} \hat{\mathbf{u}}_A + \|\mathbf{u}_{A^c}\|_1 + 2R\mathbf{u}_{A^c}^\top \boldsymbol{\Lambda}_{A^c A} \boldsymbol{\beta}_A^*.$$

Knowing that  $\hat{\mathbf{u}}_{A^c} = \mathbf{0}$  and plugging in the expression for  $\hat{\mathbf{u}}_A$ , the KKT conditions imply that

$$|-\mathbf{C}_{A^c A} \mathbf{C}_A^{-1} (\mathbf{s}_A + 2R\boldsymbol{\Lambda}_A \boldsymbol{\beta}_A^*) + 2R\boldsymbol{\Lambda}_{A^c A} \boldsymbol{\beta}_A^*| \leq 1.$$

## A.7 Proof of Theorem 4

Define

$$\begin{aligned} Z_n(\mathbf{u}) &= \sum_{i=1}^n (\epsilon_i - \mathbf{u}^\top \mathbf{x}_i / \sqrt{n})^2 - \epsilon_i^2 \\ &\quad + \lambda_1^n \sum_{j=1}^p \omega_j (|\beta_j^* + u_j / \sqrt{n}| - |\beta_j^*|) \\ &\quad + \lambda_2^n [(\boldsymbol{\beta}^* + \mathbf{u} / \sqrt{n})^\top \boldsymbol{\Lambda} (\boldsymbol{\beta}^* + \mathbf{u} / \sqrt{n}) - \boldsymbol{\beta}^{*\top} \boldsymbol{\Lambda} \boldsymbol{\beta}^*], \end{aligned}$$

which is minimized at  $\sqrt{n}(\hat{\boldsymbol{\beta}}_n^{\text{adaptive}} - \boldsymbol{\beta}^*)$ . From the proof of Theorem 1, we know that the first line in  $Z_n$  converges in distribution to  $-2\mathbf{u}^\top \mathbf{w} + \mathbf{u}^\top \mathbf{C} \mathbf{u}$ ,  $\mathbf{w} \sim N(0, \sigma^2 \mathbf{C})$ . For the second line, one has to distinguish two cases.

Case 1:  $\beta_j^* \neq 0$ . Then  $(|\beta_j^* + u_j / \sqrt{n}| - |\beta_j^*|) \rightarrow u_j \text{sign}(\beta_j^*)$ . Moreover, from the definition of  $\omega_j$ , the assumptions made for  $\hat{\beta}_j^{\text{init}}$  and the continuous mapping theorem, we have  $\omega_j \xrightarrow{P} |\beta_j^*|^{-\gamma}$ . Since  $\lambda_1^n / \sqrt{n} \rightarrow 0$  by assumption, the whole term vanishes.

Case 2:  $\beta_j^* = 0$ . Then  $\sqrt{n}(|\beta_j^* + u_j / \sqrt{n}| - |\beta_j^*|) \rightarrow |u_j|$ ,  $n^{-1/2} \lambda_1^n \omega_j = n^{-1/2} \lambda_1^n r_n^\gamma |r_n \hat{\beta}_j^{\text{init}}|^{-\gamma} \rightarrow$

$\infty$ , noting that  $|r_n \hat{\beta}_j^{\text{init}}| = O_P(1)$  by assumption. Overall, if  $u_j \neq 0$ , the whole term tends to infinity as  $n \rightarrow \infty$ . For the third line in  $Z_n(\mathbf{u})$ , one obtains the limit  $2\lambda_2^0 \mathbf{u}^\top \mathbf{\Lambda} \boldsymbol{\beta}^*$ . Putting all together, we have for all  $\mathbf{u}$  that

$$Z_n(\mathbf{u}) \xrightarrow{D} Z(\mathbf{u}) = \begin{cases} -2\mathbf{u}_A^\top \mathbf{w}_A + \mathbf{u}_A^\top \mathbf{C}_A \mathbf{u}_A + 2\lambda_2^0 \mathbf{u}_A^\top \mathbf{\Lambda}_A \boldsymbol{\beta}_A^* & \text{if } \mathbf{u}_{A^c} = \mathbf{0}, \\ \infty, & \text{otherwise,} \end{cases}$$

and that  $\text{argmin } Z_n \xrightarrow{D} \text{argmin } Z = \hat{\mathbf{u}}$  by Theorem A. 1. We have  $\hat{\mathbf{u}}_{A^c} = \mathbf{0}$ , and, by differentiation

$$\hat{\mathbf{u}}_A = \mathbf{C}_A^{-1}(\mathbf{w}_A - \lambda_2^0 \mathbf{\Lambda}_A \boldsymbol{\beta}_A^*), \quad \mathbf{w}_A \sim N(0, \sigma^2 \mathbf{C}_A),$$

i.e.

$$\sqrt{n}(\hat{\boldsymbol{\beta}}_{A,n}^{\text{adaptive}} - \boldsymbol{\beta}_A^*) \xrightarrow{D} N(-\lambda_2^0 \mathbf{C}_A^{-1} \mathbf{\Lambda}_A \boldsymbol{\beta}_A^*, \mathbf{C}_A^{-1}),$$

and consequently

$$\lim_{n \rightarrow \infty} P(\exists j \in A : \hat{\beta}_{j,n}^{\text{adaptive}} = 0) = 0.$$

On the other hand, take  $j \in A^c$  and assume that  $\lim_{n \rightarrow \infty} \hat{\beta}_{j,n}^{\text{adaptive}} \neq 0$ . Then one can differentiate the adaptive structured elastic net criterion w.r.t.  $\beta_j$ , yielding the equation

$$2\mathbf{X}_{j,n}^\top (\mathbf{y}_n - \mathbf{X}_n \hat{\boldsymbol{\beta}}_n^{\text{adaptive}}) - 2\lambda_2^n l_{jj} \hat{\beta}_{j,n}^{\text{adaptive}} - 2\lambda_2^n \sum_{r \neq j} l_{jr} \hat{\beta}_{r,n}^{\text{adaptive}} = \lambda_1^n \text{sign}(\hat{\beta}_{j,n}^{\text{adaptive}}) \omega_j.$$

Dividing both sides of the previous equation by  $\sqrt{n}$ , the left hand side is  $O_P(1)$ , while the expression on the right hand side tends to infinity. Hence, the probability that the equation is fulfilled tends to zero, which contradicts the assumption that  $\hat{\beta}_{j,n}^{\text{adaptive}} \neq 0$ .

## B Analytical solution for two predictors

First, assume that  $\hat{\beta}_1, \hat{\beta}_2 \neq 0$ . From the KKT conditions, we have

$$\begin{aligned} \hat{\beta}_1 &= \frac{\overbrace{\mathbf{X}_1^\top (\mathbf{y} - \mathbf{X}_2 \hat{\beta}_2) + \lambda_2 l_{12} \hat{\beta}_2}^{\tilde{\beta}_1} - \frac{\lambda_1}{2} \text{sign}(\tilde{\beta}_1)}{1 + \lambda_2 l_{11}}, \\ \hat{\beta}_2 &= \frac{\overbrace{\mathbf{X}_2^\top (\mathbf{y} - \mathbf{X}_1 \hat{\beta}_1) + \lambda_2 l_{22} \hat{\beta}_1}^{\tilde{\beta}_2} - \frac{\lambda_1}{2} \text{sign}(\tilde{\beta}_2)}{1 + \lambda_2 l_{22}}, \end{aligned}$$

i.e.

$$\begin{aligned} \hat{\beta}_1 &= (1 + \lambda_2 l_{11})^{-1} \text{sign}(\tilde{\beta}_1) (|\tilde{\beta}_1| - \lambda_1/2) \\ &= (1 + \lambda_2 l_{11})^{-1} \text{sign}(\tilde{\beta}_1) \left( \frac{|\tilde{\beta}_1| - |\tilde{\beta}_2|}{2} - \gamma \right), \quad \gamma = \frac{|\tilde{\beta}_1| + |\tilde{\beta}_2| - \lambda_1}{2} > 0. \end{aligned}$$

Likewise, we obtain

$$\tilde{\beta}_2 = (1 + \lambda_2 l_{22})^{-1} \text{sign}(\tilde{\beta}_2) \left( \frac{|\tilde{\beta}_2| - |\tilde{\beta}_1|}{2} - \gamma \right).$$

Otherwise, assume that  $\tilde{\beta}_1 - \lambda_1/2 > 0$  and that

$$\left[ \frac{|\tilde{\beta}_2| - |\tilde{\beta}_1|}{2} + \gamma \right]_+ = 0,$$

where  $[z]_+ = \max(z, 0)$ . Then, it follows that  $\hat{\beta}_2 = 0$ . The case  $\hat{\beta}_1 = 0$  is derived analogously. Overall, this gives the following recipe for determining  $\hat{\beta}_1, \hat{\beta}_2$ : one starts with

$$\begin{bmatrix} \tilde{\beta}_1 \\ \tilde{\beta}_2 \end{bmatrix} = \begin{bmatrix} 1 + \lambda_2 l_{11} & \mathbf{X}_1^\top \mathbf{X}_2 + \lambda_2 l_{12} \\ \mathbf{X}_1^\top \mathbf{X}_2 + \lambda_2 l_{12} & 1 + \lambda_2 l_{22} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{X}_1^\top \mathbf{y} \\ \mathbf{X}_2^\top \mathbf{y} \end{bmatrix} \begin{bmatrix} 1 + \lambda_2 l_{11} & 0 \\ 0 & 1 + \lambda_2 l_{22} \end{bmatrix},$$

computes  $\gamma$ , and determines

$$\begin{aligned} \hat{\beta}_1 &= \left[ \frac{|\tilde{\beta}_1| - |\tilde{\beta}_2|}{2} + \gamma \right]_+ (1 + \lambda_2 l_{11})^{-1}, \\ \hat{\beta}_2 &= \left[ \frac{|\tilde{\beta}_2| - |\tilde{\beta}_1|}{2} + \gamma \right]_+ (1 + \lambda_2 l_{22})^{-1}. \end{aligned}$$

## References

- Besag, J. (1974). Spatial interaction and the statistical analysis of lattice systems (with discussion). *Journal of the Royal Statistical Society Series B* 36, 192–236.
- Chung, F. (1997). *Spectral Graph Theory*. AMS Publications.
- Daumer, M., K. Thaler, E. Kruis, W. Feneberg, G. Staude, and M. Scholz (2007). Steps towards a miniaturized, robust and autonomous measurement device for the long-term monitoring of patient activity: ActiBelt. *Biomed Tech* 52, 149–155.
- Efron, B., T. Hastie, I. Johnstone, and R. Tibshirani (2004). Least Angle Regression (with discussion). *The Annals of Statistics* 32, 407–499.
- Efron, B. and R. Tibshirani (2007). On testing the significance of sets of genes. *The Annals of Applied Statistics* 1, 107–129.
- Eilers, P. and B. Marx (1996). Flexible smoothing with B-Splines and penalties (with discussion). *Statistical Science* 11, 89–121.
- Eilers, P. and B. Marx (1999). Generalized linear regression on sampled signals and curves: a P-spline approach. *Technometrics* 41, 1–13.
- Fan, J. and R. Li (2001). Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American Statistical Association* 97, 210–221.
- Frank, I. and J. Friedman (1993). A statistical view of some chemometrics regression tools (with discussion). *Technometrics* 35, 109–148.
- Friedman, J., T. Hastie, H. Hoefling, and R. Tibshirani (2007). Pathwise coordinate optimization. *The Annals of Applied Statistics* 2, 302–332.
- Genkin, A., D. Lewis, and D. Madigan (2007). Large-scale Bayesian logistic regression for text categorization. *Technometrics* 49, 589–616.



- Geyer, C. (1996). On convex stochastic optimization. Technical report, Department of Statistics, University of Minnesota.
- Goeman, J. (2007). An efficient algorithm for  $\ell^1$ -penalized estimation. Technical report, Department of Medical Statistics and Bioinformatics, University of Leiden.
- Golub, T., D. Slonim, P. Tamayo, C. Huard, M. Gaasenbeek, J. Mesirov, H. Coller, M. Loh, J. Downing, M. Caligiuri, C. Bloomfield, and E. Lander (1999). Molecular classification of cancer: class discovery and class prediction by gene expression monitoring. *Science* 286, 531–537.
- Hastie, T., A. Buja, and R. Tibshirani (1995). Penalized Discriminant Analysis. *The Annals of Statistics* 23, 73–102.
- Hoerl, A. and R. Kennard (1970). Ridge regression: biased estimation for nonorthogonal problems. *Technometrics* 8, 27–51.
- Kanehisa, M. and S. Goto (2000). KEGG: Kyoto encyclopedia of genes and genomes. *Nucleic Acids Research* 28, 27–30.
- Knight, K. and W. Fu (2000). Asymptotics for lasso-type estimators. *The Annals of Statistics* 28, 1356–1378.
- Li, C. and H. Li (2008). Network-constrained regularization and variable selection for analysis of genomic data. *Bioinformatics* 24, 1175–1182.
- McCullagh, P. and J. Nelder (1989). *Generalized Linear Models*. Chapman and Hall, London.
- Park, T. and G. Casella (2008). The Bayesian Lasso. *Journal of the American Statistical Association* 103, 681–686.
- Rosset, S., J. Zhu, and T. Hastie (2004). Boosting as a regularized path to a maximum margin classifier. *Journal of Machine Learning Research* 5, 941–973.
- Rue, H. and L. Held (2001). *Gaussian Markov Random Fields*. Chapman and Hall/CRC, Boca Raton.
- Stein, C. (1981). Estimation of the mean of a multivariate normal distribution. *The Annals of Statistics* 9, 1135–1151.
- Le Cun, Y., B. Boser, J. Denker, D. Henderson, R. Howard, W. Hubbard, and L. Jackel (1989). Backpropagation applied to handwritten zip code recognition. *Neural Computation* 2, 541–551.
- Tibshirani, R. (1996). Regression shrinkage and variable selection via the lasso. *Journal of the Royal Statistical Society Series B* 58, 671–686.
- Tibshirani, R., T. Hastie, B. Narasimhan, and G. Chu. (2003). Class prediction by nearest shrunken centroids, with applications to DNA microarrays. *Statistical Science* 18, 104–117.
- Tibshirani, R., M. Saunders, S. Rosset, J. Zhu, and K. Knight (2005). Sparsity and smoothness via the fused lasso. *Journal of the Royal Statistical Society Series B* 67, 91–108.

- Tseng, P. (2001). Convergence of a block coordinate descent method for nondifferentiable minimization. *Journal of Optimization Theory and Applications* 109, 475–494.
- Tutz, G. and J. Gertheiss (2009). Feature extraction in signal regression: a boosting technique for functional data regression. *Journal of Computational and Graphical Statistics*, to appear.
- Wang, L., J. Zhu, and H. Zhou (2006). The doubly regularized support vector machine. *Statistica Sinica* 16, 589–616.
- Wu, T. and K. Lange (2008). Coordinate descent procedures for lasso penalized regression. *The Annals of Applied Statistics* 1, 224–244.
- Zhao, P. and B. Yu (2006). On model selection consistency of the lasso. *Journal of Machine Learning Research* 7, 2541–2567.
- Zhou, H. (2006). The adaptive lasso and its oracle properties. *Journal of the American Statistical Association* 101, 1418–1429.
- Zhou, H. and T. Hastie (2005). Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society Series B* 67, 301–320.